

Kernel Methods

Prof. Nicholas Zabaras

Center for Informatics and Computational Science

<https://cics.nd.edu/>

University of Notre Dame

Notre Dame, Indiana, USA

Email: nzabaras@gmail.com

URL: <https://www.zabaras.com/>

March 22, 2019

Contents

- ❑ Keeping or discarding the training data, [Kernel Function](#), [Kernel Methods and Kernel Trick](#), [Dual Representation](#), [Constructing Kernels from Basis Functions](#), [Constructing Kernels Directly](#)
- ❑ [Kernel Engineering](#), [Combining Kernels](#), [Gaussian Kernel](#), [Kernels for Comparing Documents](#), [Matern Kernel](#), [String Kernels](#), [Symbolic Input](#), [Pyramid Match Kernels](#), [Probabilistic Generative Models](#), [Probabilistic Product Kernels](#), [Mercer Kernels](#), [Fisher Kernel](#), [Sigmoidal Kernel](#), [Radial Basis Functions](#)
- ❑ [Kernel Machines](#), [L1VM and RVM](#), [Sparse Vector Machines](#)
- ❑ [Kernel Density Estimation](#), [Smoothing Kernels for Generative Modeling](#), [KNN Classifiers](#)
- ❑ [Nadaraya-Watson model](#), [Interpolation with Noisy Inputs](#), [Kernel Regression](#), [Locally Weighted Regression](#), [K-Medoids Clustering](#)
- ❑ [Kernel PCA](#)

Following closely:

Bishop CM, [Pattern Recognition and Machine Learning](#), Springer, 2006 (Chapter 6)
K. Murphy, [Machine Learning: A probabilistic Approach](#) (Chapter 14)

Training Data: Keep or Discard?

- We have considered linear parametric methods e.g. for regression.
- The form of the mapping $y(\mathbf{x}, \mathbf{w})$ from input \mathbf{x} to output y is governed by \mathbf{w} .
- During the learning phase, the training data is used either
 - to obtain a point estimate of \mathbf{w} or
 - posterior distribution $p(\mathbf{w}|\mathbf{x}, \mathbf{t})$.
- The training data is then discarded, and predictions for new inputs are based purely on the learned parameter vector \mathbf{w} .
- The same approach is used in nonlinear parametric models including neural networks.

Training Data: Keep or Discard?

- There are pattern recognition techniques, in which the training data points are kept and used for prediction.
 - Parzen probability density model: set of kernel functions centered on training data points.
 - Nearest neighbors technique: assigning to each new test vector the same label as the closest example from the training set.
- These are **memory-based methods**: storing the training set in order to make predictions for future data points.
- They require a metric to measure the similarity of any two vectors in input space.

Kernel Functions

Kernel Function

- We define a **kernel function** to be a real-valued function of two arguments, $k(x, x') \in \mathbb{R}$, for $x, x' \in \mathcal{X}$.
- Typically the function is **symmetric** (i.e., $k(x, x') = k(x', x)$), and **non-negative** (i.e., $k(x, x') \geq 0$), so it can be interpreted as a measure of similarity (but this is not needed).
- Kernel methods are generally **fast to train but slow at making predictions** for test data points.

Kernel Methods and Kernel Trick

- Many linear parametric models can be recast into an equivalent ‘dual representation’
 - **Prediction is based on linear combinations of a kernel function evaluated at the training data points**
- For models which are based on a fixed nonlinear feature space mapping $\phi(x)$, the kernel function is given by

$$k(x, x') = \phi(x)^T \phi(x')$$

- Using the identity mapping for the feature space we obtain the **linear kernel**:

$$\phi(x) = x, \quad k(x, x') = x^T x'$$

- **Kernel trick:** if we have an algorithm formulated such that x enters only in the form of scalar products, then we can replace that scalar product with some other choice of kernel.

- Aizerman, M. A., E. M. Braverman, and L. I. Rozonoer (1964). [The probability problem of pattern recognition learning and the method of potential functions](#). *Automation and Remote Control* **25**, 1175–1190.
- Boser, B. E., I. M. Guyon, and V. N. Vapnik (1992). [A training algorithm for optimal margin classifiers](#). In D. Haussler (Ed.), *Proceedings Fifth Annual Workshop on Computational Learning Theory* (COLT), pp. 144–152. ACM.Scholkopf et al., 1998).

Types of Kernel Functions

- For models based on feature space mapping $\phi(x)$:

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

- Its a symmetric function:

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$$

- Linear kernel:

$$\phi(\mathbf{x}) = \mathbf{x}, \quad k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- Stationary kernel: invariant to translations in \mathbf{x}

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$$

- Homogeneous kernel (radial basis functions):

$$k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$$

Kernel Substitution Applications

- Algorithms expressed in terms of scalar products can be reformulated using kernel substitution (kernel trick).
Applications include:

- Non-linear PCA,
- Nearest-neighbor classifiers,
- Kernel Fisher discriminant, etc.

- Schölkopf, B., A. Smola, and K.-R. Müller (1998). [Nonlinear component analysis as a kernel eigenvalue problem](#). *Neural Computation* **10**(5), 1299–1319.
- Mika, S., G. Rätsch, J. Weston, and B. Schölkopf (1999). [Fisher discriminant analysis with kernels](#). In Y. H. Hu, J. Larsen, E. Wilson, and S. Douglas (Eds.), [Neural Networks for Signal Processing IX](#), pp. 41–48. IEEE.
- Roth, V. and V. Steinlage (2000). [Nonlinear discriminant analysis using kernel functions](#). In S. A. Solla, T. K. Leen, and K. R. Müller (Eds.), [Advances in Neural Information Processing Systems](#), Volume 12. MIT Press.
- Baudat, G. and F. Anouar (2000). [Generalized discriminant analysis using a kernel approach](#). *Neural Computation* **12**(10), 2385–2404.

Textbooks on Kernel Methods include:

- Schölkopf, B. and A. J. Smola (2002). [Learning with Kernels](#). MIT Press.
- Herbrich, R. (2002). [Learning Kernel Classifiers](#). MIT Press.
- Shawe-Taylor, J. and N. Cristianini (2004). [Kernel Methods for Pattern Analysis](#). Cambridge University Press ([slides](#))

Dual Representation

Dual Representation

- Many linear models for regression and classification can be reformulated in terms of a dual representation in which the kernel function arises naturally.
- Consider a linear regression model with a regularized cost functional ($\lambda > 0$).

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Set the gradient of the first term to zero:

$$\begin{aligned}\mathbf{w} &= -\frac{1}{\lambda} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a} \\ a_n &\equiv -\frac{1}{\lambda} (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n), \quad \mathbf{a} = \{a_1, \dots, a_N\}^T\end{aligned}$$

- Note that we defined \mathbf{a} using only the 1st term in $J(\mathbf{w})$.

Dual Representation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}, \quad \mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) - t_n) \boldsymbol{\phi}(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

□ Substitute \mathbf{w} in the cost functional:

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2} \sum_{n=1}^N (\mathbf{a}^T \Phi \boldsymbol{\phi}(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T \Phi \sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) t_n + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \end{aligned}$$

□ Here we used from an earlier lecture

$$\Phi = \begin{pmatrix} \boldsymbol{\phi}^T(\mathbf{x}_1) \\ \boldsymbol{\phi}^T(\mathbf{x}_2) \\ \vdots \\ \boldsymbol{\phi}^T(\mathbf{x}_N) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\phi}_0(\mathbf{x}_1) & \boldsymbol{\phi}_1(\mathbf{x}_1) & .. & \boldsymbol{\phi}_{M-1}(\mathbf{x}_1) \\ \boldsymbol{\phi}_0(\mathbf{x}_2) & \boldsymbol{\phi}_1(\mathbf{x}_2) & .. & \boldsymbol{\phi}_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \vdots & \vdots \\ \boldsymbol{\phi}_0(\mathbf{x}_N) & \boldsymbol{\phi}_1(\mathbf{x}_N) & .. & \boldsymbol{\phi}_{M-1}(\mathbf{x}_N) \end{pmatrix}, \quad \begin{aligned} \Phi^T \Phi &= \sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T, \\ \Phi^T \mathbf{t} &= \sum_{n=1}^N t_n \boldsymbol{\phi}(\mathbf{x}_n) \end{aligned}$$

Dual Representation: Gramm Matrix

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

- We can now define the **Gramm matrix (symmetric)** $\mathbf{K}_{N \times N}$ using the kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$ as follows:

$$K_{nm} = K_{mn} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) \equiv k(\mathbf{x}_n, \mathbf{x}_m) \quad \text{or} \quad \mathbf{K} = \Phi \Phi^T$$

$$\mathbf{K} = \Phi \Phi^T = \begin{pmatrix} \phi^T(\mathbf{x}_1) \\ \phi^T(\mathbf{x}_2) \\ \vdots \\ \phi^T(\mathbf{x}_N) \end{pmatrix} (\phi(\mathbf{x}_1) \quad \phi(\mathbf{x}_2) \quad \dots \quad \phi(\mathbf{x}_N))$$

- The cost function takes the **dual representation** form:

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

Solution of the Dual Representation

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

- Setting the gradient wrt \mathbf{a} equal to zero:

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_{N \times N})^{-1} \mathbf{t}, \quad K_{nm} = K_{mn} = k(\mathbf{x}_n, \mathbf{x}_m) \quad \text{or} \quad \mathbf{K} = \Phi \Phi^T \quad (N \times N \text{ matrix})$$

- Using $\mathbf{w} = \Phi^T \mathbf{a}$, the prediction for a new input \mathbf{x} is

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \phi(\mathbf{x})^T \Phi^T \mathbf{a} = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_{N \times N})^{-1} \mathbf{t}$$

where

$$\begin{aligned} \mathbf{k}(\mathbf{x})^T &= \phi(\mathbf{x})^T \Phi^T = (\phi(\mathbf{x})^T \phi(\mathbf{x}_1) \quad \phi(\mathbf{x})^T \phi(\mathbf{x}_2) \quad \dots \quad \phi(\mathbf{x})^T \phi(\mathbf{x}_N)) \\ &= (k(\mathbf{x}_1, \mathbf{x}) \quad k(\mathbf{x}_2, \mathbf{x}) \quad \dots \quad k(\mathbf{x}_N, \mathbf{x})) \end{aligned}$$

- We now invert an $N \times N$ matrix (to determine \mathbf{a}) instead of $M \times M$ (to compute \mathbf{w}). Everything is in terms of kernels (many choices here, kernel trick) but often $N \gg M$.

Summary of the Dual Representation

- The dual formulation allows the solution to the problem to be expressed entirely in terms of the kernel function $k(\mathbf{x}, \mathbf{x}')$.
- The prediction at \mathbf{x} is given by a linear combination of the target values from the training set.

$$y(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_{N \times N})^{-1} \mathbf{t}$$
$$\mathbf{k}(\mathbf{x})^T = (k(\mathbf{x}_1, \mathbf{x}) \quad k(\mathbf{x}_2, \mathbf{x}) \quad \dots \quad k(\mathbf{x}_N, \mathbf{x}))$$

- Because we now have to invert a larger matrix $N \times N$, the formulation does not seem practical.
- We now have the choice of working with kernels rather than the basis functions $\phi(\mathbf{x})$.

- ✓ The feature space can be infinite dimensional.
- The cost of computing \mathbf{a} is $\mathcal{O}(N^3)$, whereas the cost of computing \mathbf{w} is $\mathcal{O}(M^3)$. However, prediction using the dual variables takes $\mathcal{O}(N)$ time, while prediction with the primal variables only takes $\mathcal{O}(M)$ time.

Dual of the Dual Representation

- It is not difficult to see that the dual of the dual formulation is given by our original least squares representation for \mathbf{w} .
- Note that $J(\mathbf{a})$ can be written in terms of $\mathbf{K}\mathbf{a}$, $\mathbf{K} = \Phi\Phi^T$

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2}\mathbf{a}^T\Phi\Phi^T\Phi\Phi^T\mathbf{a} - \mathbf{a}^T\Phi\Phi^T\mathbf{t} + \frac{1}{2}\mathbf{t}^T\mathbf{t} + \frac{\lambda}{2}\mathbf{a}^T\Phi\Phi^T\mathbf{a} = \\ &= \frac{1}{2}(\mathbf{K}\mathbf{a} - \mathbf{t})^T(\mathbf{K}\mathbf{a} - \mathbf{t}) + \frac{\lambda}{2}\mathbf{a}^T\mathbf{K}\mathbf{a} \end{aligned}$$

- Recall that:

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} \phi^T(x_1) \\ \phi^T(x_2) \\ \vdots \\ \phi^T(x_N) \end{pmatrix} (\phi(x_1) \quad \phi(x_2) \quad .. \quad \phi(x_N))$$

- Since $N \gg M$, \mathbf{K} is rank deficient – $N - M$ eigenvectors of \mathbf{K} have zero eigenvalues. We can decompose \mathbf{a} as: $\mathbf{a} = \mathbf{a}_{\parallel} + \mathbf{a}_{\perp}$, $\mathbf{a}_{\perp}^T\mathbf{a}_{\parallel} = 0$ and $\mathbf{K}\mathbf{a}_{\perp} = \mathbf{0}$. Thus \mathbf{a}_{\perp} is not determined by $J(\mathbf{a})$. One can set $\mathbf{a}_{\perp} = \mathbf{0}$, and thus

$$\mathbf{a} = \mathbf{a}_{\parallel} \in \text{span } \mathbf{K} \quad (\mathbf{K} = \Phi\Phi^T) \Rightarrow \mathbf{a} = \Phi(\Phi^T \mathbf{v}) \Rightarrow \mathbf{a} = \Phi \mathbf{u}$$

Dual of the Dual Representation

- The cost functional becomes:

$$J(\mathbf{u}) = \frac{1}{2} (\Phi \Phi^T \Phi \mathbf{u} - \mathbf{t})^T (\Phi \Phi^T \Phi \mathbf{u} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{u}^T \Phi^T \Phi \Phi^T \Phi \mathbf{u}$$

- Since $\Phi^T \Phi$ has full rank, we can introduce an equivalent parametrization as

$$\mathbf{w} = \Phi^T \Phi \mathbf{u}$$

- The cost functional now becomes:

$$J(\mathbf{u}) = \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w}$$

- This is our original least squares problem. So the proof is complete.

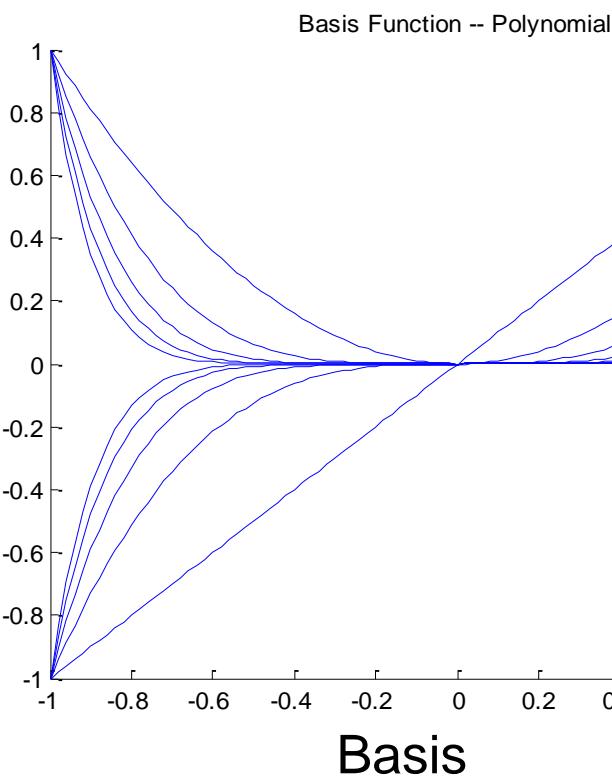
Constructing Kernels from Basis Functions

Constructing Kernels From Basis Functions

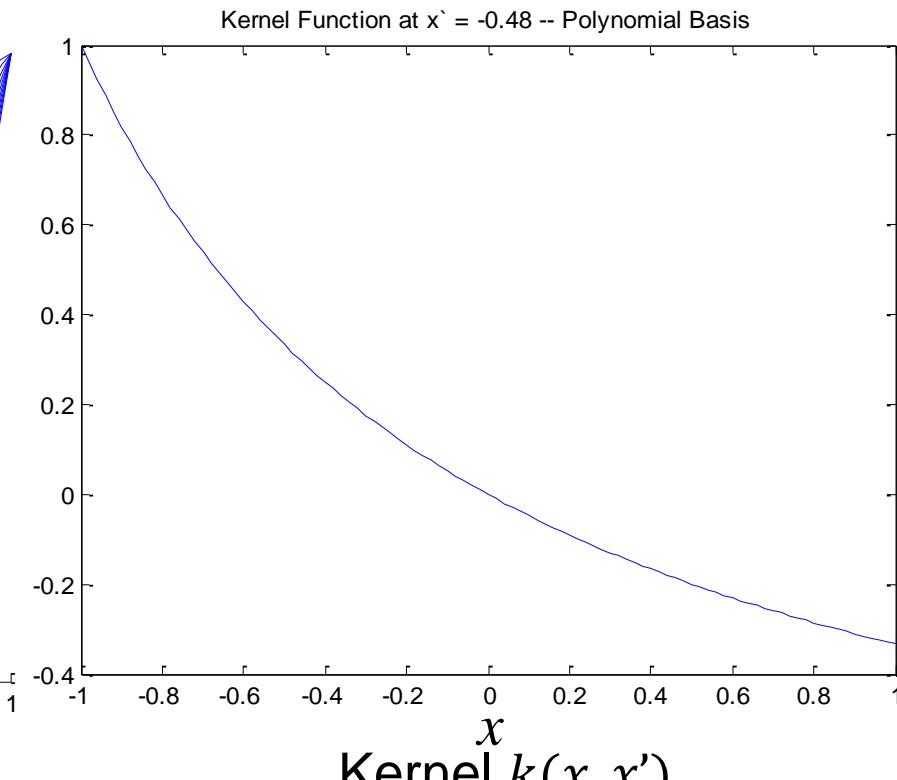
- Choose feature space mapping $\phi(x)$, then the kernel (in 1D) is given as:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

Polynomials



MatLab Code

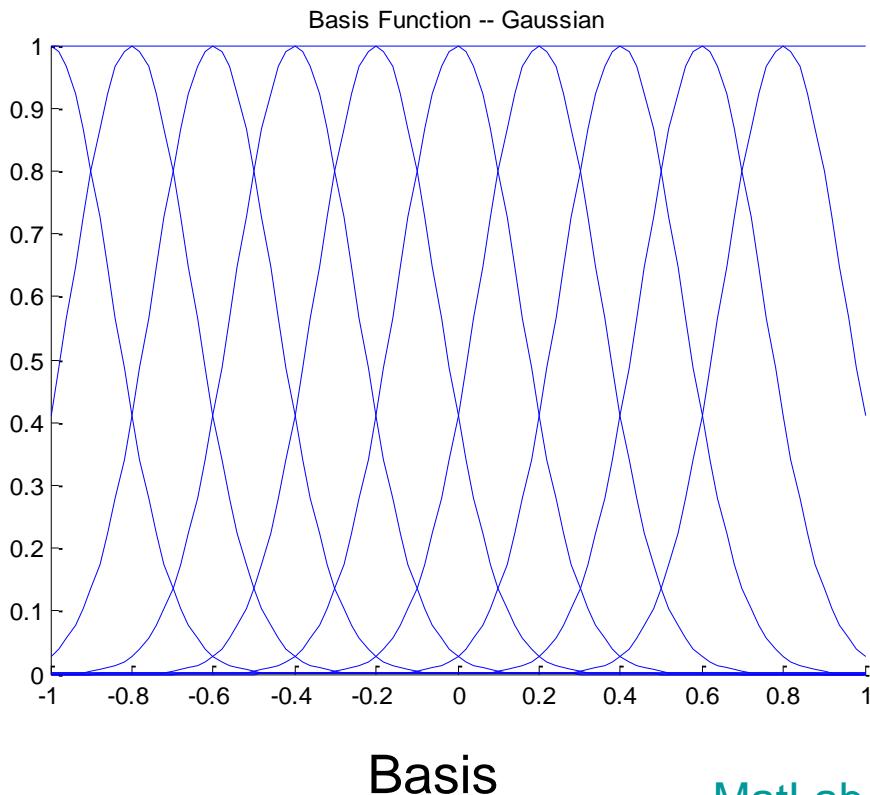


Constructing Kernels From Basis Functions

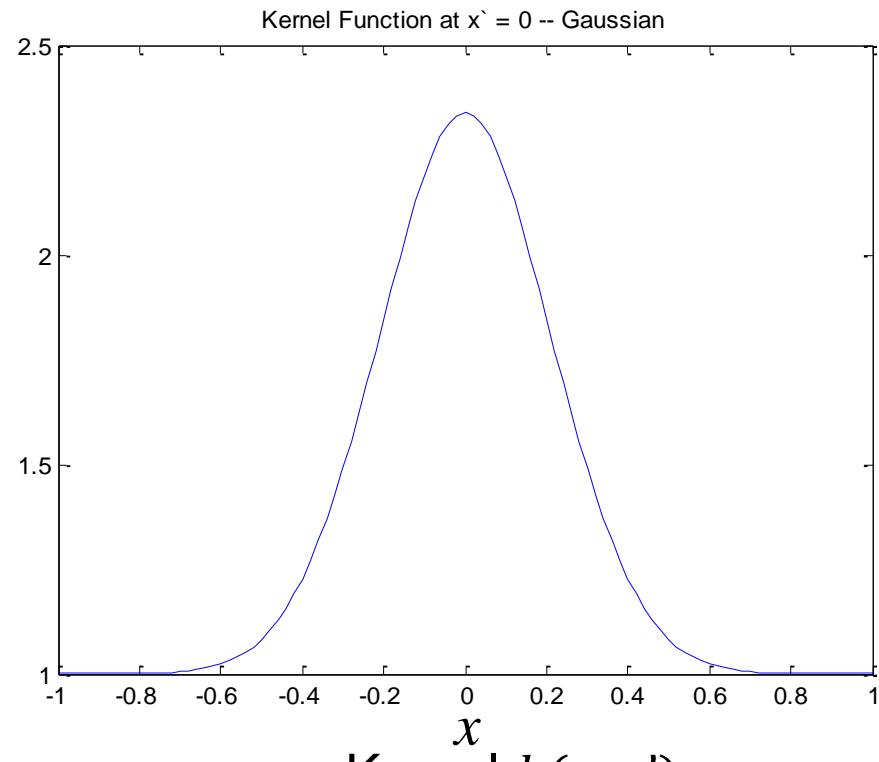
- Choose feature space mapping $\phi(x)$, then the kernel (in 1D) is given as:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

Gaussians



MatLab Code

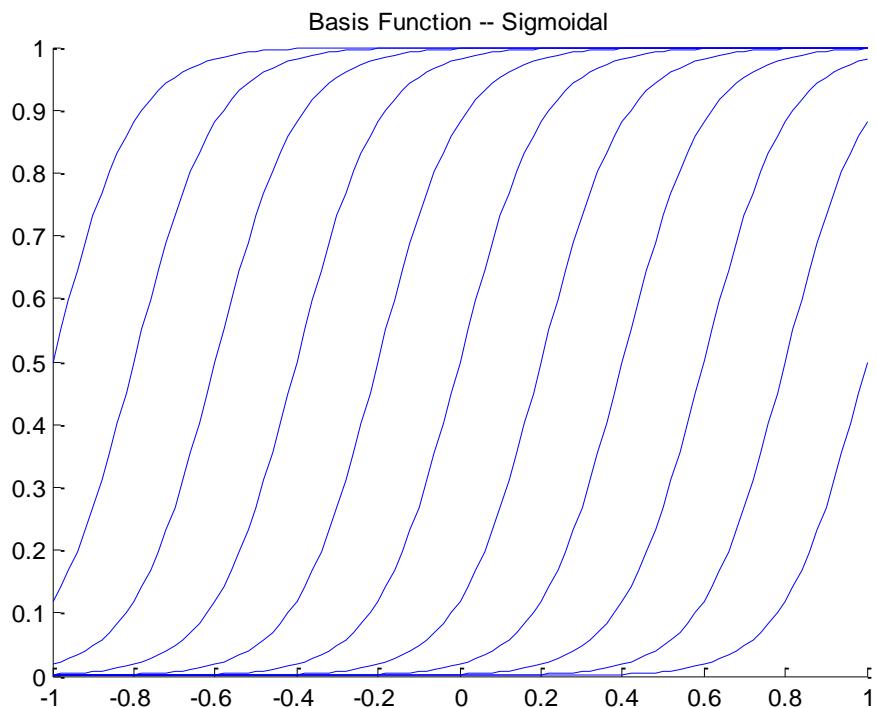


Constructing Kernels From Basis Functions

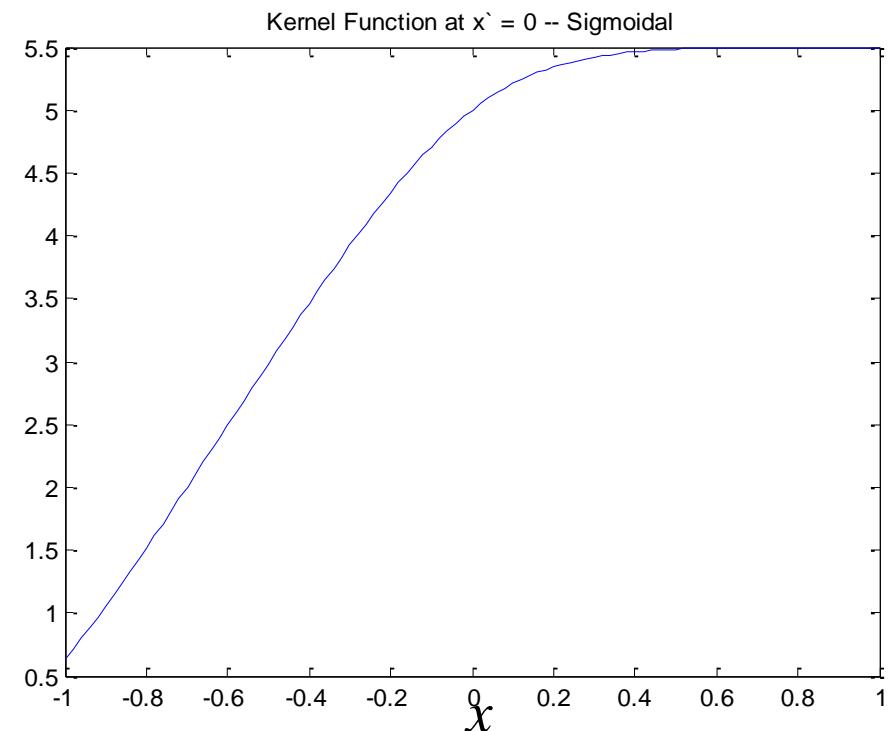
- Choose feature space mapping $\phi(x)$, then the kernel (in 1D) is given as:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

Logistic Sigmoid



MatLab Code



Constructing Kernels Directly

- When constructing the kernel functional directly verify its validity
 - it should correspond to a scalar product in some feature space

- Simple example

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

- In the 2-D case this corresponds to

$$k(\mathbf{x}, \mathbf{z}) = (x_1 z_1 + x_2 z_2)^2 = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

where

$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^T$$

- To test validity without having to construct $\phi(\mathbf{x})$ explicitly, one can check if $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$ is semidefinite:

Function $k(\mathbf{x}_n, \mathbf{x}_m)$ is a valid kernel $\Leftrightarrow K \geq 0$ (Gram matrix) for all $\{\mathbf{x}_n\}$

- Shawe-Taylor, J. and N. Cristianini (2004). *Kernel Methods for Pattern Analysis*. Cambridge University Press ([slides](#))

Kernel Engineering

Kernel Engineering

- Build kernels using simple kernels as building blocks. Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$, $k_2(\mathbf{x}, \mathbf{x}')$, the following kernels are valid:

$$\begin{array}{ll} k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}'), & k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \\ k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')), & k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \\ k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}'), & k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \\ k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')), & k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T A \mathbf{x}' \\ k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b), & k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \end{array}$$

with corresponding conditions:

$c > 0$, $f(\cdot)$ arbitrary, $q(\cdot)$ polynomial with ≥ 0 coefficients,

$\phi(\mathbf{x}) \in \mathbb{R}^M$, $k_3(\cdot, \cdot)$ valid kernel in \mathbb{R}^M , A sym. ≥ 0 ,

$\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, k_a , k_b valid kernels

Combining Kernels: Kernel Engineering

- If $k_1(x, x')$ is a valid kernel, so is $k(x, x') = ck_1(x, x')$, $c > 0$

$$k(x, x') = ck_1(x, x') = c\phi(x)^T \phi(x') = u(x)^T u(x'), \text{ where } u(x) = \sqrt{c}\phi(x)$$

- If $k_1(x, x')$ is a valid kernel so is $k(x, x') = f(x)k_1(x, x')f(x')$ for $f(\cdot)$ arbitrary.

$$k(x, x') = f(x)k_1(x, x')f(x') = f(x)\phi(x)^T \phi(x')f(x') = u(x)^T u(x'), \text{ where } u(x) = f(x)\phi(x)$$

- If $k_1(x, x'), k_2(x, x')$ are valid kernels so is $k_1(x, x') + k_2(x, x')$
Indeed, note that the Gram matrix K with elements $k(x_n, x_m)$ for this kernel is ≥ 0 .

$$\forall a, a^T Ka = a^T K_1 a + a^T K_2 a \geq 0, \text{ since } K_1, K_2 \geq 0$$

- Schoelkopf, B. and A. Smola (2002). [Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond](#). MIT Press.

Kernel Examples

□ Polynomial kernels:

- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$: $\phi(\mathbf{x})$ contains only terms of degree 2
- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^2, c > 0$: $\phi(\mathbf{x})$ contains constant, linear and order 2 terms
- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^M$: contains all monomials of order M

Example: For \mathbf{x} and \mathbf{x}' images, $k(\mathbf{x}, \mathbf{x}')$ is a weighted sum of all possible products of M pixels of image \mathbf{x} with M pixels of image \mathbf{x}' .

- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M, c > 0$: contains all terms up to order M

Squared Exponential or Gaussian Kernel

- The **squared exponential kernel** (SE kernel) or **Gaussian kernel** is defined by

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \Sigma^{-1} (\mathbf{x} - \mathbf{x}')\right)$$

- If Σ is diagonal, this can be written as

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2} \sum_{j=1}^D \frac{1}{\sigma_j^2} (x_j - x'_j)^2\right)$$

- σ_j defines the **characteristic length scale** of dimension j . If $\sigma_j = \infty$, the j -th dimension is ignored; this is known as the **ARD kernel**.
- If Σ is spherical, we get the isotropic kernel $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$
- Here σ^2 is known as the **bandwidth** and the kernel as the **radial basis function** or **RBF kernel**, since it is only a function of $\|\mathbf{x} - \mathbf{x}'\|$.

RBF or Gaussian Kernel

- Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

- To see that this is a valid kernel note:

$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}'$$

and write it as:

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}')\right) \\ &= \underbrace{\exp\left(-\frac{\mathbf{x}^T \mathbf{x}}{2\sigma^2}\right)}_{f(\mathbf{x})} \underbrace{\exp\left(\frac{\mathbf{x}^T \mathbf{x}'}{\sigma^2}\right)}_{k_1(\mathbf{x}, \mathbf{x}')} \underbrace{\exp\left(-\frac{\mathbf{x}'^T \mathbf{x}'}{2\sigma^2}\right)}_{f(\mathbf{x}')} \end{aligned}$$

- Expanding the middle exponential, we can see that the Gaussian kernel can be expressed as the inner product of an infinite dimensional feature vector.

Gaussian Kernel: non-Euclidean Distance

- Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}')\right)$$

- With [the kernel trick](#), we can replace the Euclidean distance with any non-linear kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2}(\kappa(\mathbf{x}, \mathbf{x}) + \kappa(\mathbf{x}', \mathbf{x}') - 2\kappa(\mathbf{x}, \mathbf{x}'))\right)$$

Kernels for Comparing Documents

- Consider documents x_i and $x_{i'}$ using a bag of words representation, x_{ij} is the number of times word j occurs in document i .
- One can use a **cosine similarity** (cosine of the angle between $x_i, x_{i'}^T$):

$$k(x_i, x_{i'}) = \frac{x_i^T x_{i'}}{\|x_i\|_2 \|x_{i'}\|_2}$$

- For orthogonal vectors (no common words) $k(x_i, x_{i'}) = 0$. In general

$$0 \leq k(x_i, x_{i'}) \leq 1$$

- This measure unfortunately accounts for document similarity that may be due to the appearance of non discriminative common words (**stop words**) in documents such as 'the' or 'and'.
- In addition, if a discriminative word occurs many times in a document, the similarity is artificially boosted (once a word is used in a document it is very likely to be used again).

Kernels for Comparing Documents

- To improve performance we replace the word count x_{ij} with a new feature vector (TF-IDF) - “term frequency - inverse document frequency”. Define the TF as:

$$tf(x_{ij}) \equiv \log(1 + x_{ij})$$

- This reduces the impact of words that occur many times within one document. The inverse document frequency is defined as

$$idf(j) \equiv \log\left(\frac{N}{1 + \sum_{i=1}^N \mathbb{I}(x_{ij} > 0)}\right)$$

- N is the total number of documents, and the denominator counts how many documents contain term j . Finally, we define the feature space

$$tf-idf(j) \equiv \phi(x_i) \equiv [tf(x_{ij}) \times idf(j)]_{j=1}^V$$

- The kernel is then taken as: $k(x_i, x_{i'}) = \frac{\phi(x_i)^T \phi(x_{i'})}{\|\phi(x_i)\|_2 \|\phi(x_{i'})\|_2}$

- Elkan, C. (2005). [Deriving TF-IDF as a Fisher kernel](#). In *Proc. Intl. Symp. on String Processing and Information Retrieval (SPIRE)*, pp. 296–301.
- Manning, C. & H. Schuetze (1999) [Foundations of statistical natural language processing](#) MIT Press.

Matern Kernel

- This is mostly used in Gaussian Processes. It takes the form:

$$k(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\ell} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}r}{\ell} \right)$$

- Here K_ν is a modified Bessel function, $\nu > 0, \ell > 0, r = \|x - x'\|$
- As $\nu \rightarrow \infty$, the Matern kernel approaches the SE kernel.
- For $\nu = \frac{1}{2}$, $k(r) = \exp(-r/\ell)$

If $D = 1$, the corresponding GP with this kernel is equivalent to the Ornstein-Uhlenbeck process (Brownian motion, continuous but not differentiable)

String Kernels

- s is a **substring** of x if we can write $x = usv$ for some strings u, s, v .
- Let $\phi_s(x)$ denote the number of times that substring s appears in string x . We define the kernel between two strings x and x' as

$$\kappa(x, x') = \sum_{s \in \mathcal{A}^*} w_s \phi_s(x) \phi_s(x'), w_s \geq 0$$

and \mathcal{A}^* is the set of all strings (of any length) from the alphabet \mathcal{A} . This is a Mercer kernel, and can be computed in $\mathcal{O}(|x| + |x'|)$ time (for certain settings of the weights $\{w_s\}$) using suffix trees.

- For $w_s = 0$ for $|s| > 1$, we get a **bag-of-characters kernel**. $\phi_1(x)$ is number of times each character in \mathcal{A} occurs in x .
- If s is to be bordered by white-space, we get a bag-of-words kernel, where $\phi_s(x)$ counts how many times each possible word occurs.

String Kernels

$$\kappa(x, x') = \sum_{s \in \mathcal{A}^*} w_s \phi_s(x) \phi_s(x'), w_s \geq 0$$

- If we consider strings of a fixed length k , we get the **k-spectrum kernel**.
- Can allow character mismatches, string kernels to compare trees, etc.

- Leslie, C., E. Eskin, A. Cohen, J. Weston, and W. Noble (2003). [Mismatch string kernels for discriminative protein classification](#). *Bioinformatics* 1, 1–10.
- Vishwanathan, S. V. N. and A. Smola (2003). [Fast kernels for string and tree matching](#). In *NIPS*.
- Shawe-Taylor, J. and N. Cristianini (2004). *Kernel Methods for Pattern Analysis*. Cambridge.
- Collins, M. and N. Duffy (2002). [Convolution kernels for natural language](#). In *NIPS*.

Kernels for Symbolic Inputs

- The inputs can be symbolic, graphs, sets, strings, text, etc.
- For example, fix a set D and consider all of its subsets ($2^{|D|}$). If A_1 and A_2 are two subsets of D , we can define the kernel as:

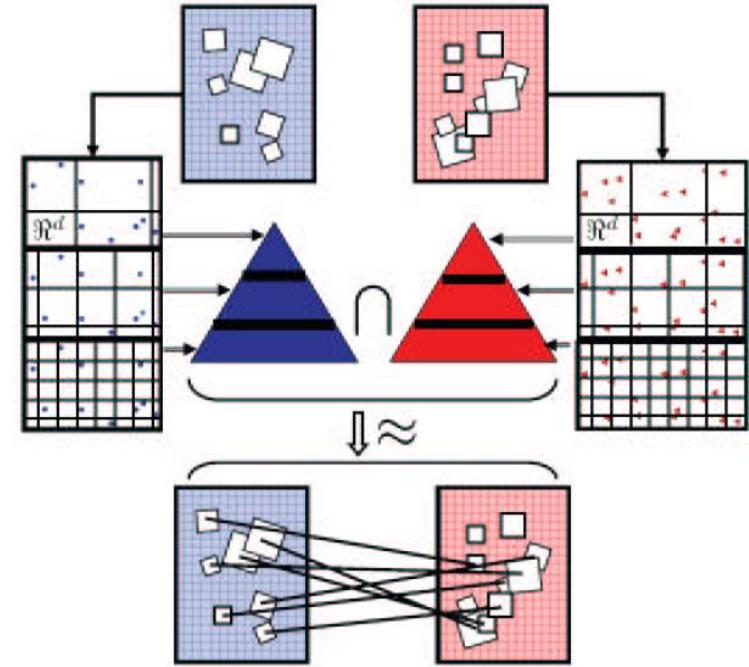
$$k(A_1, A_2) = 2^{|A_1 \cap A_2|}$$

where $|A_1 \cap A_2|$ is the number of elements in the intersection $A_1 \cap A_2$

- $\phi(D)$ is defined to map to a vector of $2^{|D|}$ 1's, one for each possible subset of D , including D itself as well as the empty set.
- For $A \subset D$, $\phi(A)$ will have 1's in all positions that correspond to subsets of A and 0's in all other positions. Therefore, $\phi(A_1)^T \phi(A_2)$ will count the number of subsets shared by A_1 and A_2 .
- This can also be obtained by counting the number of elements in the intersection of A_1 and A_2 , and then raising 2 to this number, which is exactly $k(A_1, A_2) = 2^{|A_1 \cap A_2|}$

Pyramid Match Kernels

- Each feature set is mapped to a multi-resolution histogram.
- These are then compared using weighted histogram intersection.
- This provides a good approximation to the similarity measure one would obtain by performing an optimal bipartite match at the finest spatial resolution, and then summing up pairwise similarities between matched points.
- The histogram method is faster and is more robust to missing and unequal numbers of points.



- Lowe, D. G. (1999). [Object recognition from local scale-invariant features](#). In *Proc. of the International Conference on Computer Vision ICCV*, Corfu, pp. 1150–1157.
- Grauman, K. and T. Darrell (2007, April). [The Pyramid Match Kernel: Efficient Learning with Sets of Features](#). *J. of Machine Learning Research* 8, 725–760.

Probabilistic Generative Models

- Use a generative model to define a kernel and then use this kernel in a discriminative approach.
- Kernel for generative model $p(\mathbf{x})$:

$$k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}')$$

- This is a valid kernel function since it is an inner product in the 1D feature space defined by the mapping $p(\mathbf{x})$.
- It says that two inputs \mathbf{x} and \mathbf{x}' are similar if they both have high probabilities.

- Haussler, D. (1999). [Convolution kernels on discrete structures](#). Technical Report UCSC-CRL-99-10, University of California, Santa Cruz, Computer Science Department.
- Lasserre, J., C. M. Bishop, and T. Minka (2006). [Principled hybrids of generative and discriminative models](#). In [Proceedings 2006 IEEE Conference on Computer Vision and Pattern Recognition](#), New York.
- Jebara, T., R. Kondor, and A. Howard (2004). [Probability product kernels](#). *J. of Machine Learning Research* 5, 819–844.

Probabilistic Generative Models

- Using $k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}')$ and the kernel engineering formulas

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}'), c > 0$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

we can generalize as:

$$k(\mathbf{x}, \mathbf{x}') = \sum_i p(\mathbf{x} | i)p(\mathbf{x}' | i)p(i)$$

$$k(\mathbf{x}, \mathbf{x}') = \int p(\mathbf{x} | z)p(\mathbf{x}' | z)p(z)dz$$

i and z are latent variables.

- \mathbf{x} and \mathbf{x}' will give large value of $k(\mathbf{x}, \mathbf{x}')$ if they have significant probability under a range of components.

Probabilistic Generative Models

- Now consider an observation is given by a sequence $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_L\}$.
- A popular generative model for sequences is the [Hidden Markov model \(HMM\)](#), which expresses the distribution $p(\mathbf{X})$ as a marginalization over a corresponding sequence of hidden states $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_L\}$.
- We define [a kernel function measuring the similarity of two sequences \$\mathbf{X}\$ and \$\mathbf{X}'\$](#) by extending the earlier mixture representation

$$k(\mathbf{X}, \mathbf{X}') = \sum_{\mathbf{Z}} p(\mathbf{X} | \mathbf{Z}) p(\mathbf{X}' | \mathbf{Z}) p(\mathbf{Z})$$

- Note that [both observed sequences \$\mathbf{X}, \mathbf{X}'\$ are generated by the same hidden sequence \$\mathbf{Z}\$](#) .

Probabilistic Product Kernels

- Jebara et al. defined the so called probability product kernel as follows:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \int p(\mathbf{x} / \mathbf{x}_i)^\rho p(\mathbf{x} / \mathbf{x}_j)^\rho d\mathbf{x}, \rho > 0$$

- $p(\mathbf{x} | \mathbf{x}_i)$ is often approximated by $p(\mathbf{x} | \widehat{\boldsymbol{\theta}}(\mathbf{x}_i))$, where $\widehat{\boldsymbol{\theta}}(\mathbf{x}_i)$ is a parameter estimate computed using a single data vector.
- Note that the fitted model to a single data point is only being used to see how similar two objects are. If we fit the model to \mathbf{x}_i and then the model thinks \mathbf{x}_j is likely, this means that \mathbf{x}_i and \mathbf{x}_j are similar.
- For example, let $p(\mathbf{x} | \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$, where σ^2 is constant. If $\rho = 1$, and we use $\widehat{\boldsymbol{\mu}}(\mathbf{x}_i) = \mathbf{x}_i$ and $\widehat{\boldsymbol{\mu}}(\mathbf{x}_j) = \mathbf{x}_j$, then we obtain an RBF kernel:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{(4\pi\sigma^2)^{D/2}} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{4\sigma^2}\right)$$

- Jebara, T., R. Kondor, and A. Howard (2004). [Probability product kernels](#). *J. of Machine Learning Research* 5, 819–844.

Mercer (Positive Definite) Kernels

- We define a Mercer (or positive definite) kernel as one for which the Gram matrix K is positive definite for all inputs $\mathbf{x}_i, i = 1, \dots, N$

$$K = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

- Both the Gaussian and cosine similarity kernels are Mercer kernels.

- Schoelkopf, B. and A. Smola (2002). [Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond](#). MIT Press.
- Sahami, M. and T. Heilman (2006). [A Web-based Kernel Function for Measuring the Similarity of Short Text Snippets](#). In *WWW conference*.

Mercer (Positive Definite) Kernels

- **Mercer's theorem.** If the Gram matrix is positive definite, we can compute an eigenvector decomposition of it as follows

$$\mathbf{K} = \mathbf{U}^T \Lambda \mathbf{U}$$

where Λ is a diagonal matrix of eigenvalues $\lambda_i > 0$. We can write:

$$K_{ij} = (\Lambda^T \mathbf{U}_{:i})^T (\Lambda^T \mathbf{U}_{:j}) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j), \text{ where } \phi(\mathbf{x}_i) \equiv \Lambda^T \mathbf{U}_{:i}$$

- Thus we see that the entries in the kernel matrix can be computed by performing an inner product of some **feature vectors** that are implicitly defined by the eigenvectors \mathbf{U} .
- If the kernel is Mercer, then there exists a function ϕ mapping $\mathbf{x} \in \mathcal{X}$ to \mathbb{R}^M such that $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$, where ϕ depends on the eigen *functions* of k (M is a potentially infinite dimensional space).

Mercer Kernels

- Consider the (non-stationary) polynomial kernel $k(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^T \mathbf{x}' + r)^M$ where $r > 0$. The corresponding feature vector $\phi(\mathbf{x})$ contains all terms up to degree M . E.g., for $M = 2$, $\gamma = r = 1$ and $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^2$, we have

$$(1 + \mathbf{x}^T \mathbf{x}')^2 =$$

$$(1 + x_1 x'_1 + x_2 x'_2)^2 = 1 + 2x_1 x'_1 + 2x_2 x'_2 + (x_1 x'_1)^2 + (x_2 x'_2)^2 + 2x_1 x'_1 x_2 x'_2$$

- Thus $\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2)^T$ i.e. this kernel is equivalent to working in a 6-dimensional feature space.

- In the case of a Gaussian kernel, the feature map lives in an infinite dimensional space.
- An example of a kernel that is not a Mercer kernel (\mathbf{K} is not positive semi-definite) is the so-called **sigmoid kernel**, defined by

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\gamma \mathbf{x}^T \mathbf{x}' + r)$$

- Can build up new Mercer kernels from simpler ones using a set of standard rules. E.g., if k_1 and k_2 are both Mercer, so is $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$.

Fisher Kernel

- Consider a parametric generative model $p(\mathbf{x}|\boldsymbol{\theta})$. We want a kernel that measures the similarity of two input vectors \mathbf{x} and \mathbf{x}' induced by this distribution.
- The Fisher score is introduced defining a vector in future space

$$g(\boldsymbol{\theta}, \mathbf{x}) = \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x} | \boldsymbol{\theta})$$

- The Fisher kernel is defined as

$$k(\mathbf{x}, \mathbf{x}') = g(\boldsymbol{\theta}, \mathbf{x})^T \mathbf{F}^{-1} g(\boldsymbol{\theta}, \mathbf{x}')$$

The Fisher information matrix \mathbf{F} (considers the geometry of the parameter/future space) is given as the expectation under $p(\mathbf{x} | \boldsymbol{\theta})$

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} \left[g(\boldsymbol{\theta}, \mathbf{x}) g(\boldsymbol{\theta}, \mathbf{x})^T \right]$$

- Using \mathbf{F} makes the kernel invariant under $\boldsymbol{\theta} \rightarrow \psi(\boldsymbol{\theta})$ where $\psi(\boldsymbol{\theta})$ is non-linear reparametrization that is invertible & differentiable.

- Jaakkola, T. S. and D. Haussler (1999). [Exploiting generative models in discriminative classifiers](#). In M. S. Kearns, S. A. Solla, and D. A. Cohn (Eds.), [Advances in Neural Information Processing Systems](#), Volume 11. MIT Press.
- Amari, S. I. (1998). [Natural gradient works efficiently in learning](#). [Neural Computation 10, 251–276](#).

Fisher Kernel

$$\mathbf{F} = \mathbb{E}_x \left[g(\theta, x) g(\theta, x)^T \right]$$

- We often approximate the information matrix with a sample average:

$$\mathbf{F} \approx \frac{1}{N} \sum_{n=1}^N g(\theta, \mathbf{x}_n) g(\theta, \mathbf{x}_n)^T$$

- This is the covariance matrix of the Fisher scores, and so the Fisher kernel corresponds to a whitening of these scores.

$$k(\mathbf{x}, \mathbf{x}') = g(\theta, \mathbf{x})^T \mathbf{F}^{-1} g(\theta, \mathbf{x}')$$

- More simply, we can just omit the Fisher information matrix altogether and use the non-invariant kernel

$$k(\mathbf{x}, \mathbf{x}') = g(\theta, \mathbf{x})^T g(\theta, \mathbf{x}')$$

▪ Hofmann, T. (2000). [Learning the similarity of documents: an information-geometric approach to document retrieval and classification](#). In S. A. Solla, T. K. Leen, and K. R. Müller (Eds.), [Advances in Neural Information Processing Systems](#), Volume 12, pp. 914–920. MIT Press.

Fisher Kernel

- The **score vector** $g(\theta, x)$ as well as the Fisher information matrix F are often computed at the MLE estimate of the parameters $\hat{\theta}$
- Note that $\hat{\theta}$ is a function of all the data, so the similarity of x and x' is computed in the context of all the data as well. Also, note that we only have to fit one model.
- The intuition behind the Fisher kernel is the following: let $g(\hat{\theta}, x)$ be the direction (in parameter space) in which x would like the parameters to move (from $\hat{\theta}$) so as to maximize its own likelihood; call this the directional gradient.
- Then we say that two vectors x and x' are similar if their directional gradients are similar wrt the geometry encoded by the curvature of the likelihood function.

Fisher Kernel

- Interestingly, it was shown that the string kernel discussed earlier is equivalent to the Fisher kernel derived from an L 'th order Markov chain.
- Also a kernel defined by the inner product of TF-IDF vectors is approximately equal to the Fisher kernel for a certain generative model of text based on the compound Dirichlet multinomial model.

- Elkan, C. (2005). [Deriving TF-IDF as a Fisher kernel](#). In *Proc. Intl. Symp. on String Processing and Information Retrieval (SPIRE)*, pp. 296–301.
- Saunders, C., J. Shawe-Taylor, and A. Vinokourov (2003). [String Kernels, Fisher Kernels and Finite State Automata](#). In *NIPS*.

Fisher Kernel: Example

- Let us compute the Fisher kernel for the Gaussian distribution with fixed covariance S :

$$p(\mathbf{x} | \boldsymbol{\mu}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, S)$$

- The Fisher score is given:

$$g(\boldsymbol{\mu}, \mathbf{x}) = \nabla_{\boldsymbol{\mu}} \ln p(\mathbf{x} | \boldsymbol{\mu}) = S^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

- The Fisher information matrix \mathbf{F} is then:

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} \left[g(\boldsymbol{\mu}, \mathbf{x}) g(\boldsymbol{\mu}, \mathbf{x})^T \right] = S^{-1} \underbrace{\mathbb{E}_{\mathbf{x}} \left[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right]}_S S^{-1} = S^{-1}$$

- The Fisher kernel is then given as:

$$k(\mathbf{x}, \mathbf{x}') = g(\boldsymbol{\mu}, \mathbf{x})^T \mathbf{F}^{-1} g(\boldsymbol{\mu}, \mathbf{x}') = (\mathbf{x} - \boldsymbol{\mu})^T S^{-1} (\mathbf{x}' - \boldsymbol{\mu})$$

- This is the Mahalanobis distance!

Sigmoidal Kernel

- Another kernel function is the sigmoidal kernel

$$k(\mathbf{x}, \mathbf{x}') = \tanh(a\mathbf{x}^T \mathbf{x}' + b)$$

- Unfortunately the Gram matrix in general is not positive semidefinite.
- This kernel was used because it gives kernel expansions such as the SVM a superficial resemblance to neural network models.
- As we shall see later in this course, in the limit of an infinite number of basis functions, a Bayesian neural network with an appropriate prior reduces to a Gaussian process, thereby providing a link between neural networks and kernel methods.

- Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer

Radial Basis Functions

- By definition

$$\phi_j(\mathbf{x}) = h(\|\mathbf{x} - \boldsymbol{\mu}_j\|)$$

where $\|\mathbf{x} - \boldsymbol{\mu}_j\|$ is a radial distance from a center $\boldsymbol{\mu}_j$.

- Originally were introduced for the problem of exact interpolation $f(\mathbf{x}_n) = t_n, n = 1, \dots, N$

$$f(\mathbf{x}) = \sum_{n=1}^N w_n h(\|\mathbf{x} - \mathbf{x}_n\|)$$

- There are as many unknowns w_n as target values t_n . Thus least squares leads to exact interpolation.
- If the target values t_n are noisy, exact interpolation overfits.

- Powell, M. J. D. (1987). [Radial basis functions for multivariable interpolation: a review](#). In J. C. Mason and M. G. Cox (Eds.), *Algorithms for Approximation*, pp. 143–167. Oxford University Press.

Radial Basis: Regularized Least Squares

- For a sum-of-squares error function with a regularizer defined in terms of a differential operator, the optimal solution is given by an expansion in the Green's functions of the operator (analogous to eigenvectors of a discrete matrix) again with one basis function centered on each data point.
- Green's functions for an isotropic differential operator as regularizer lead to radial basis functions (depending only on the distance from the data points)
- Interpolation is not exact due to the presence of the regularizer.

- Poggio, T. and F. Girosi (1990). [Networks for approximation and learning](#). *Proceedings of the IEEE* **78**(9), 1481–1497.
- Bishop, C. M. (1995a). [Neural Networks for Pattern Recognition](#). Oxford University Press

Kernel Machines

Kernel Machines

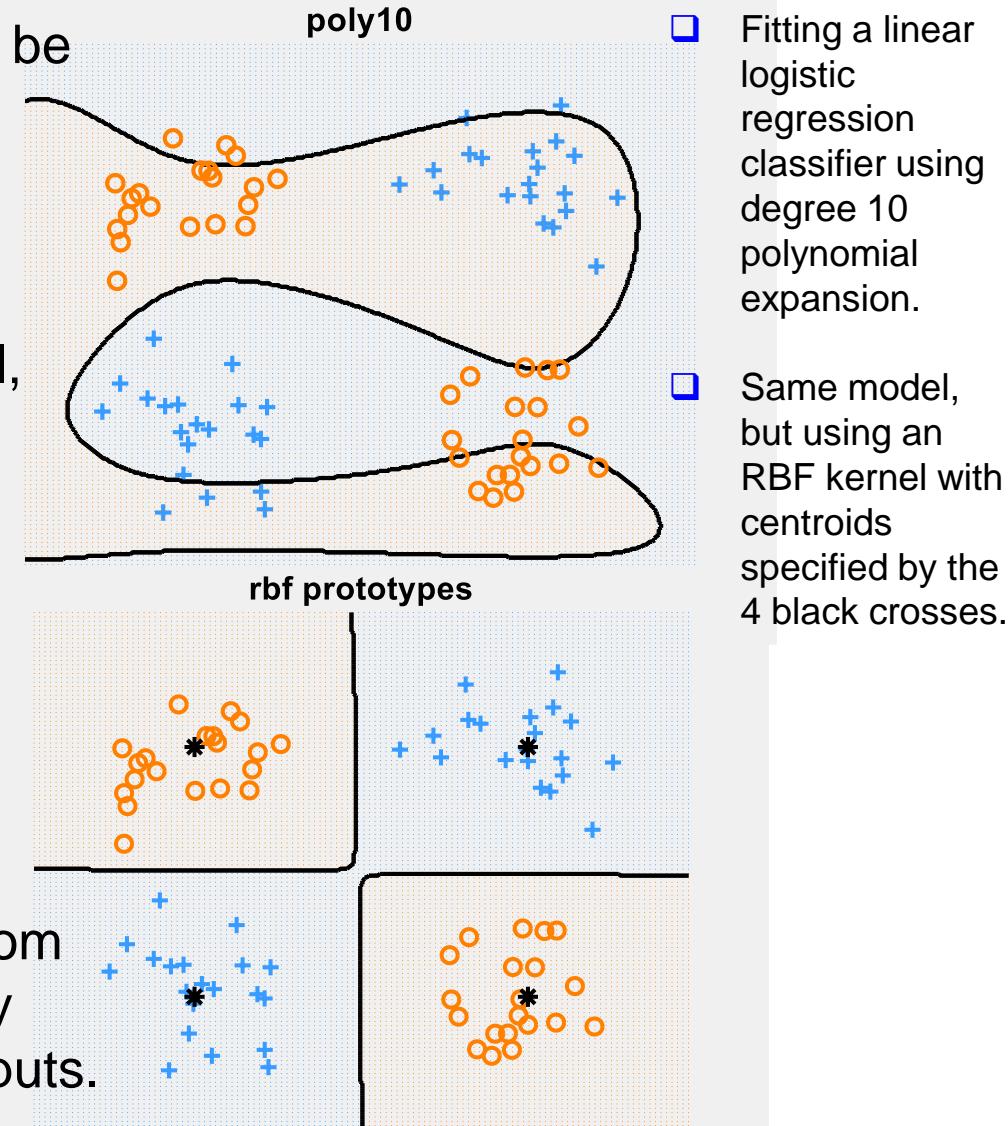
- We define a **kernel machine** to be a GLM where the input feature vector has the form

$$\phi(x) = [k(x, \mu_1), \dots, k(x, \mu_K)]$$

where $\mu_i \in X$ are a set of K **centroids**. If k is an RBF kernel, this is called an **RBF network**.

- This defines a kernelized feature vector. The kernel need not be a Mercer kernel.
- Consider the kernelized feature vector for logistic regression by defining $p(y|x, \theta) = \text{Ber}(w^T \phi(x))$.
- We consider the data coming from the **xor** function. This is a binary valued function of two binary inputs.

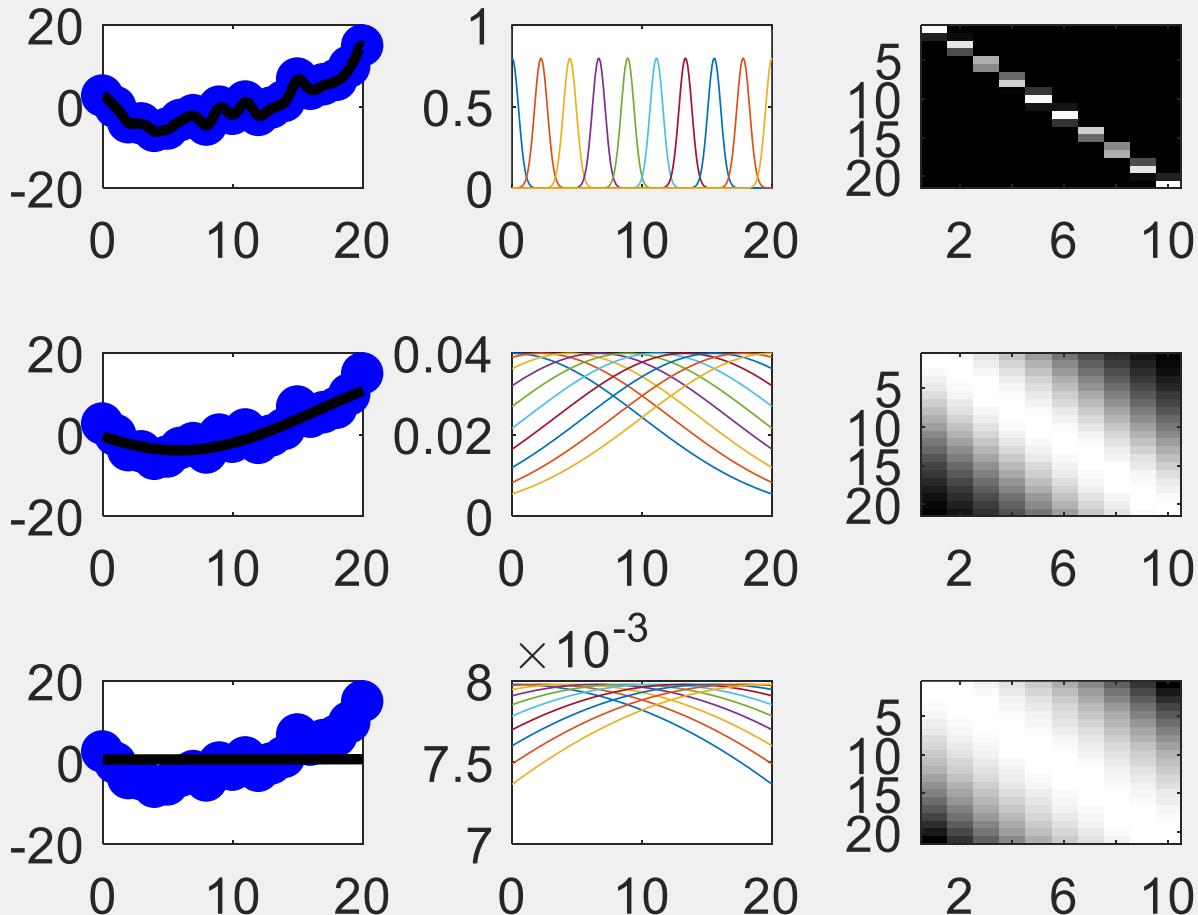
[logregXorDemo.m](#)
from [PMTK3](#)



Kernelized Feature Vector in Regression

- RBF basis in 1d.
- Left column: fitted function.
- Middle column: basis functions evaluated on a grid.
- Right column: design matrix.
- K=10 uniformly spaced RBF. Top to bottom we show different bandwidths:
 $\tau = 0.1, \tau = 0.5, \tau = 50.$

We can use the kernelized feature vector inside a linear regression model by defining $p(y|x, \theta) = \mathcal{N}(\mathbf{w}^T \phi(x), \sigma^2)$

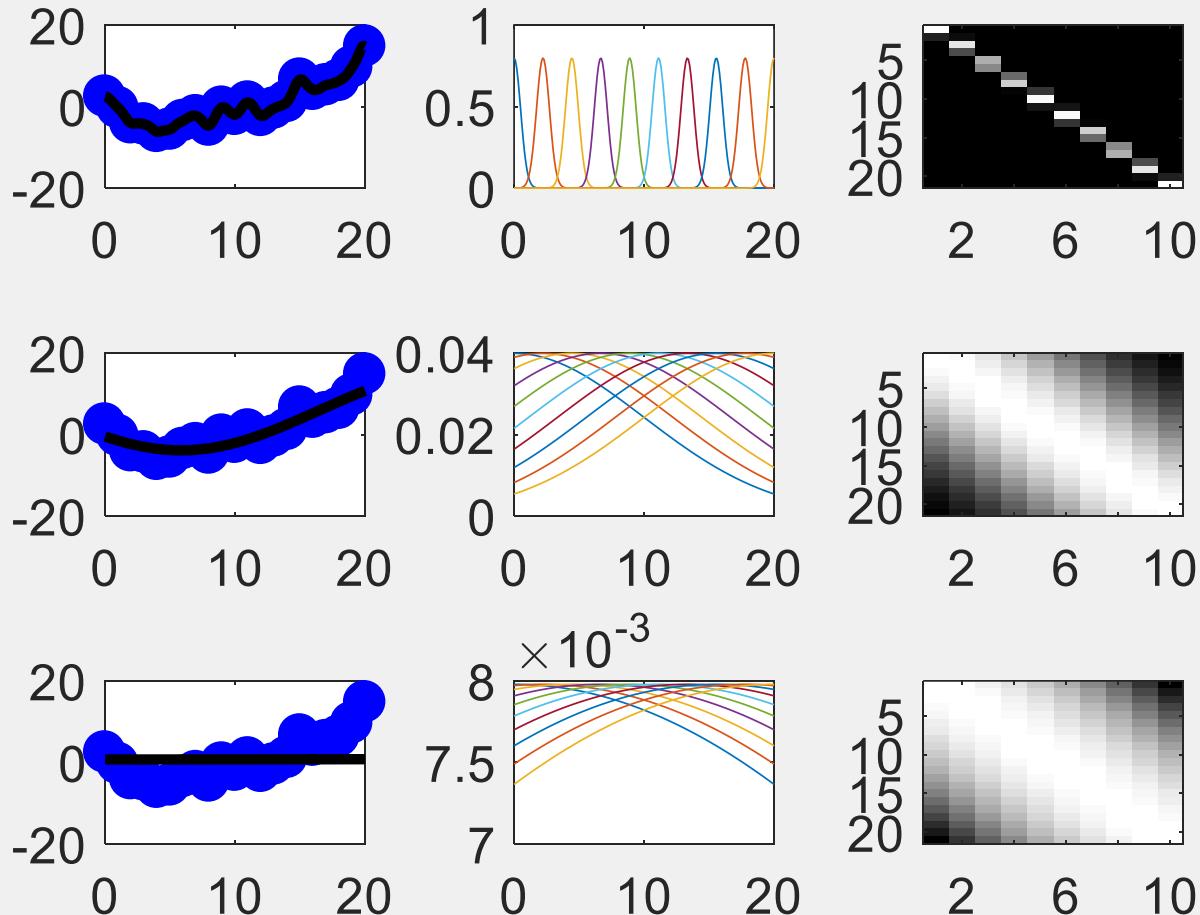


[linregRbfDemo.m](#)
from [PMTK3](#)

Kernelized Feature Vector in Regression

- Small τ leads to very wiggly functions. The predicted function value will only be non-zero for x that are close to one of the μ_k .
- If τ is very large, the design matrix reduces to a constant matrix of 1's. Each point is equally close to every prototype; the corresponding function is just a straight line.

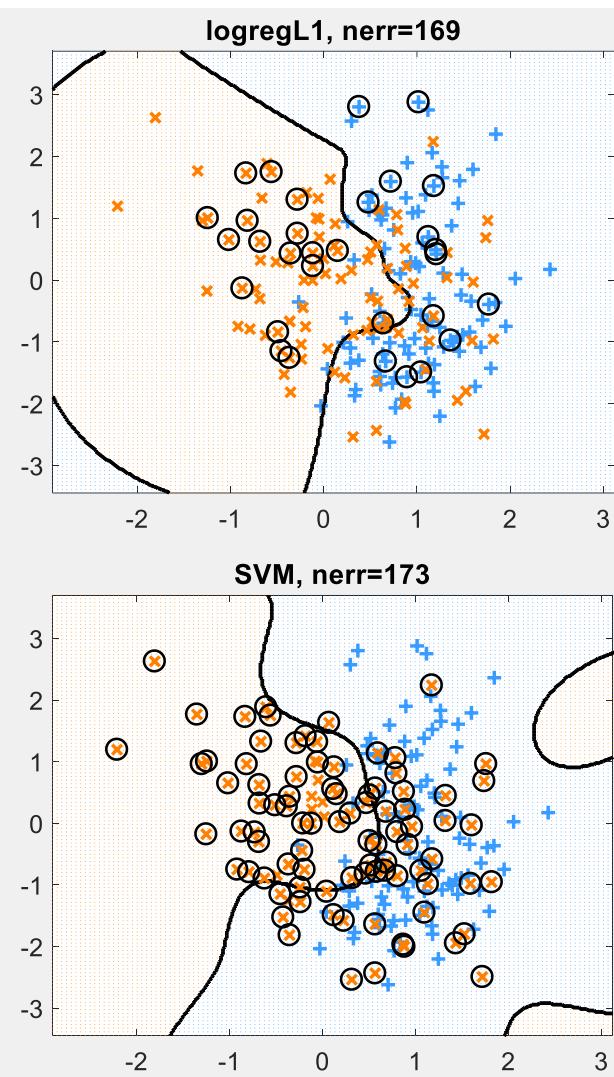
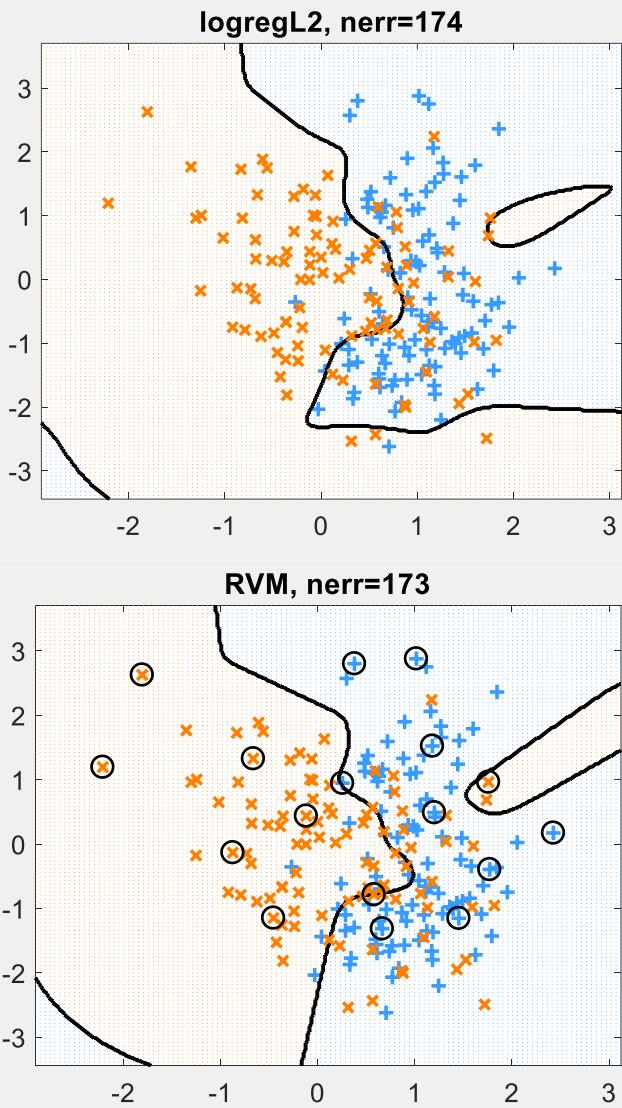
We define $p(y|x, \theta) = \mathcal{N}(\mathbf{w}^T \phi(x), \sigma^2)$.



[linregRbfDemo.m](#)
from [PMTK3](#)

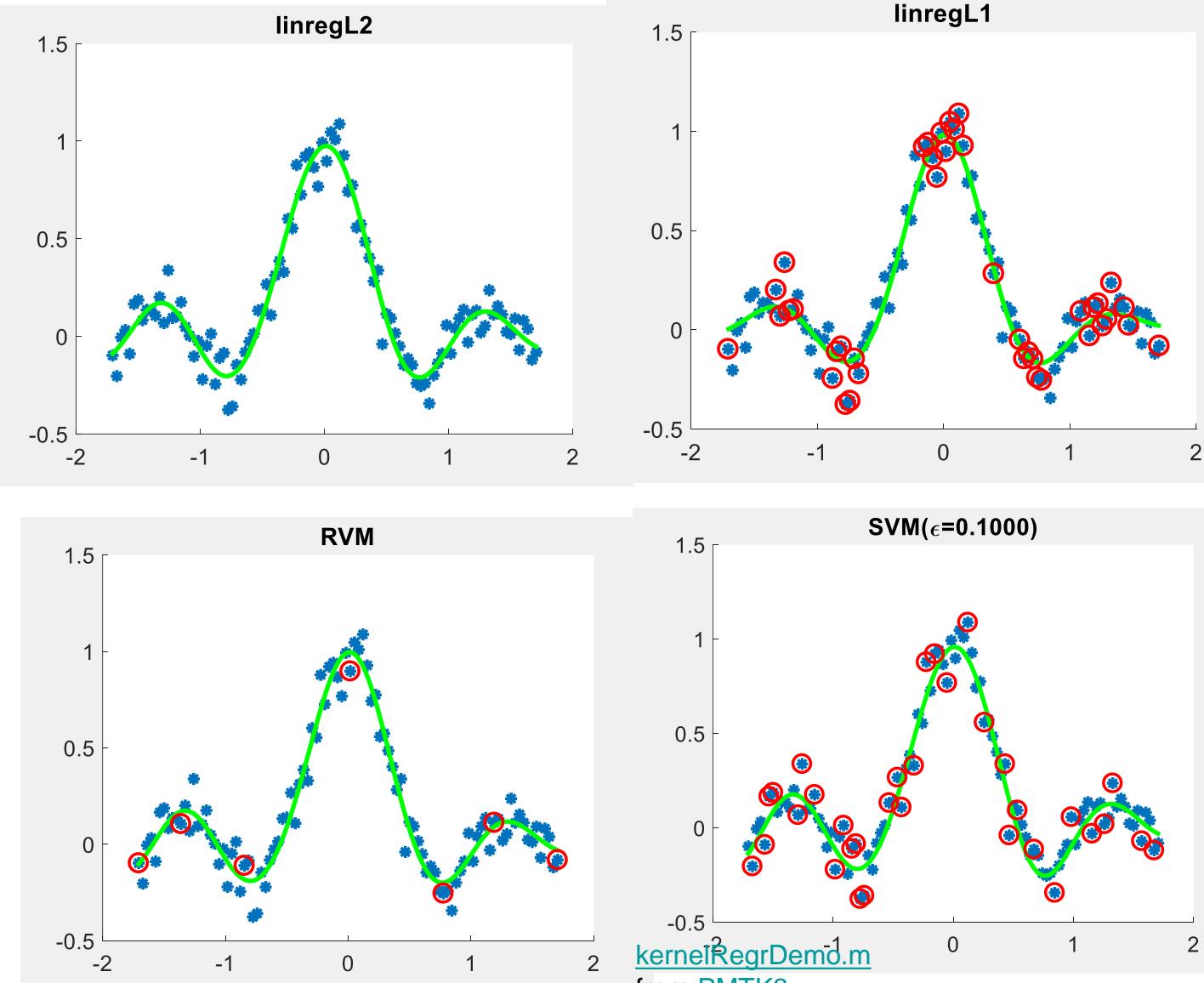
Sparse Vector Machines

- ❑ Non-linear binary classification using an RBF kernel with $\sigma = 0.3$.
- ❑ (a) L2VM with $\lambda = 5$.
- ❑ (b) L1VM with $\lambda = 1$.
- ❑ (c) RVM.
- ❑ (d) SVM with $C = 1/\lambda$ chosen by cross validation.
- ❑ Black circles denote the support vectors.
- ❑ All the methods give similar performance.
- ❑ RVM is the fastest to train and sparsest (and fastest at test time), then L1VM, and then SVM.



Sparse Vector Machines

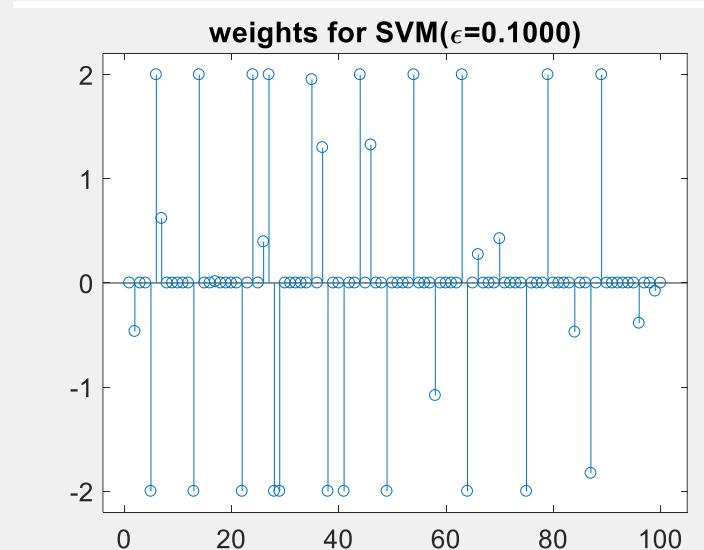
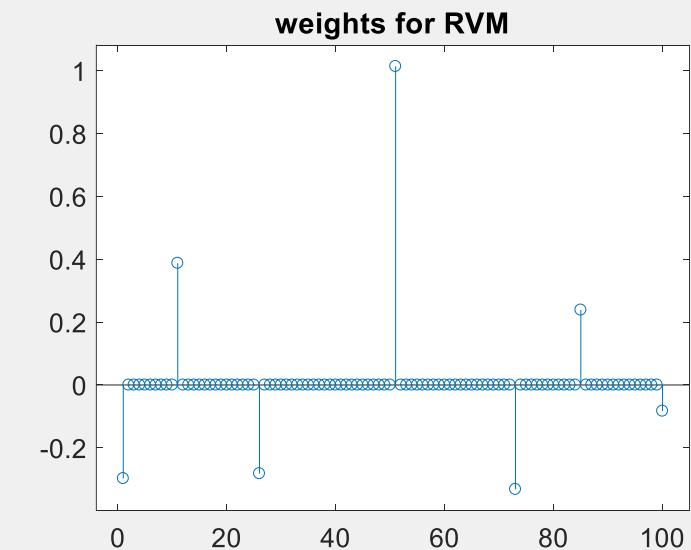
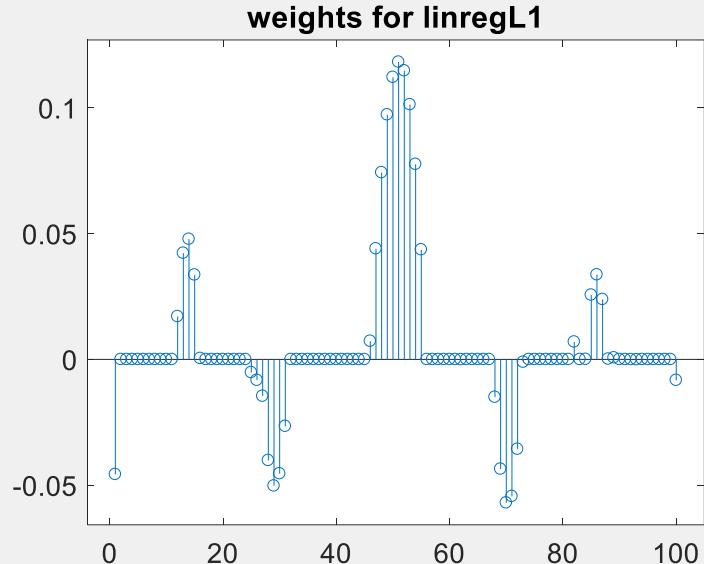
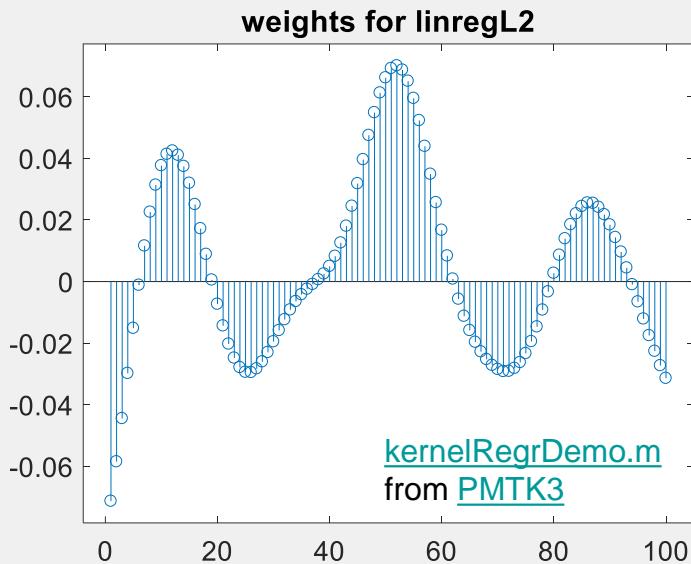
- We compare L2VM, L1VM, RVM and an SVM.
- The noisy sinc function using an RBF kernel with $\sigma = 0.3$.
- (a) L2VM with $\lambda = 0.5$.
- (b) L1VM with $\lambda = 0.5$.
- (c) RVM.
- (d) SVM with $C = 1/\lambda$ chosen by cross validation, and $\epsilon = 0.1$. Red circles denote the retained training examples.
- Predictions are similar, but RVM is the sparsest, then L2VM, then SVM.



[kernelRegrDemo.m](#)
from [PMTK3](#)

Sparse Vector Machines

- ❑ Coefficient vectors of length $N=100$



Kernel Density Estimation

Smoothing Kernels for Generative Models

- Smoothing kernels are used for non-parametric density estimates. This can be used for unsupervised density estimation, $p(x)$, as well as for generative models $p(y, x)$.
- A smoothing kernel has one argument and satisfies

$$\int k(x)dx = 1, \int xk(x)dx = 0, \int x^2k(x)dx > 0$$

- Gaussian kernel:

$$k(x) \equiv \frac{1}{(2\pi)^{1/2}} e^{-x^2/2}$$

$$k_h(x) \equiv \frac{1}{h} k\left(\frac{x}{h}\right), k_h(x) \equiv k_h(\|x\|)$$

$$k_h(x) \equiv \frac{1}{h^D (2\pi)^{D/2}} \prod_{j=1}^D \exp\left(-\frac{1}{2h^2} x_j^2\right)$$

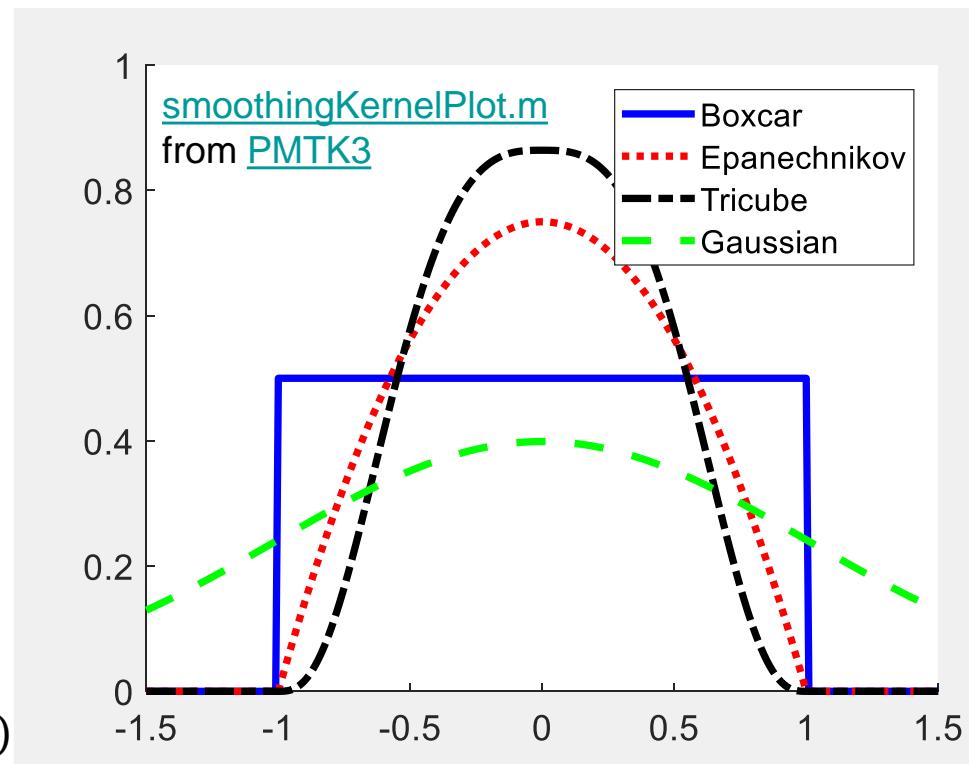
- Epanechnikov kernel:

$$k(x) \equiv \frac{3}{4} (1 - x^2) \mathbb{I}(|x| \leq 1)$$

- Tri-Cube kernel:

$$k(x) \equiv \frac{70}{81} (1 - |x|^3) \mathbb{I}(|x| \leq 1)$$

- Boxcar Kernel: $k(x) \equiv \mathbb{I}(|x| \leq 1)$



Kernel Density Estimation

- One can build a non-parametric density estimator by using a Gaussian mixture with as many components as the available data points, and with means taking to correspond to each data point:

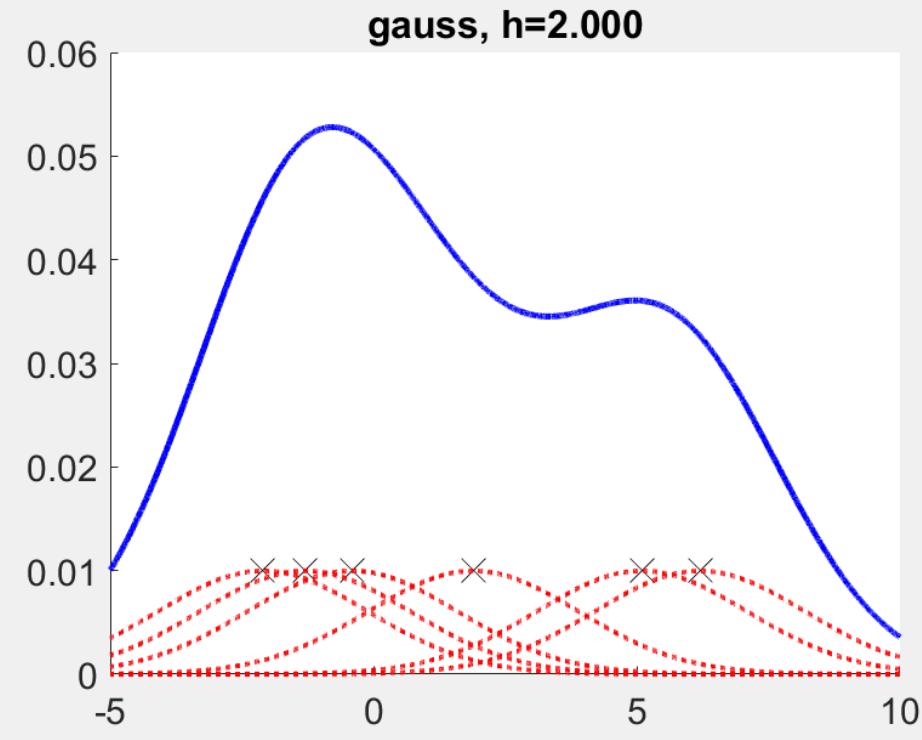
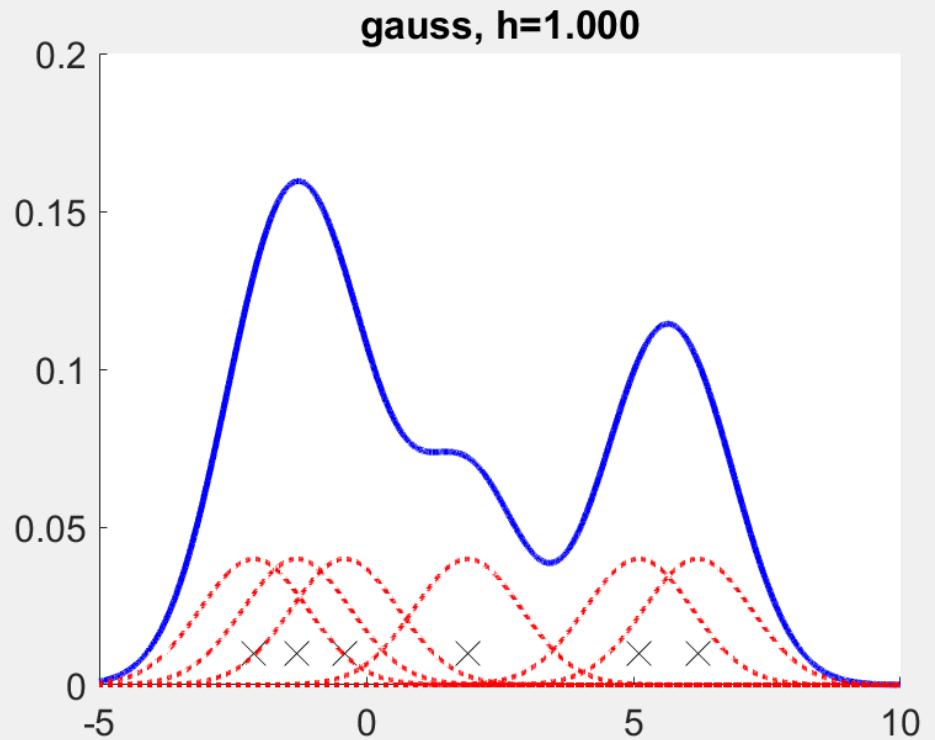
$$p(\mathbf{x}|\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \mathcal{N}(\mathbf{x}|\mathbf{x}_i, \sigma^2 \mathbf{I})$$

- This can be generalized using smooth kernels as:

$$p(\mathbf{x}|\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N k_h(\mathbf{x} - \mathbf{x}_i)$$

- This is called a [Parzen window density estimator](#), or [kernel density estimator \(KDE\)](#).
- Advantages: No model fitting is required (except for tuning the bandwidth, done by cross-validation) and no need to pick K .
- Disadvantages: Takes a lot of memory to store, and a lot of time to evaluate. It is also of no use for clustering tasks.

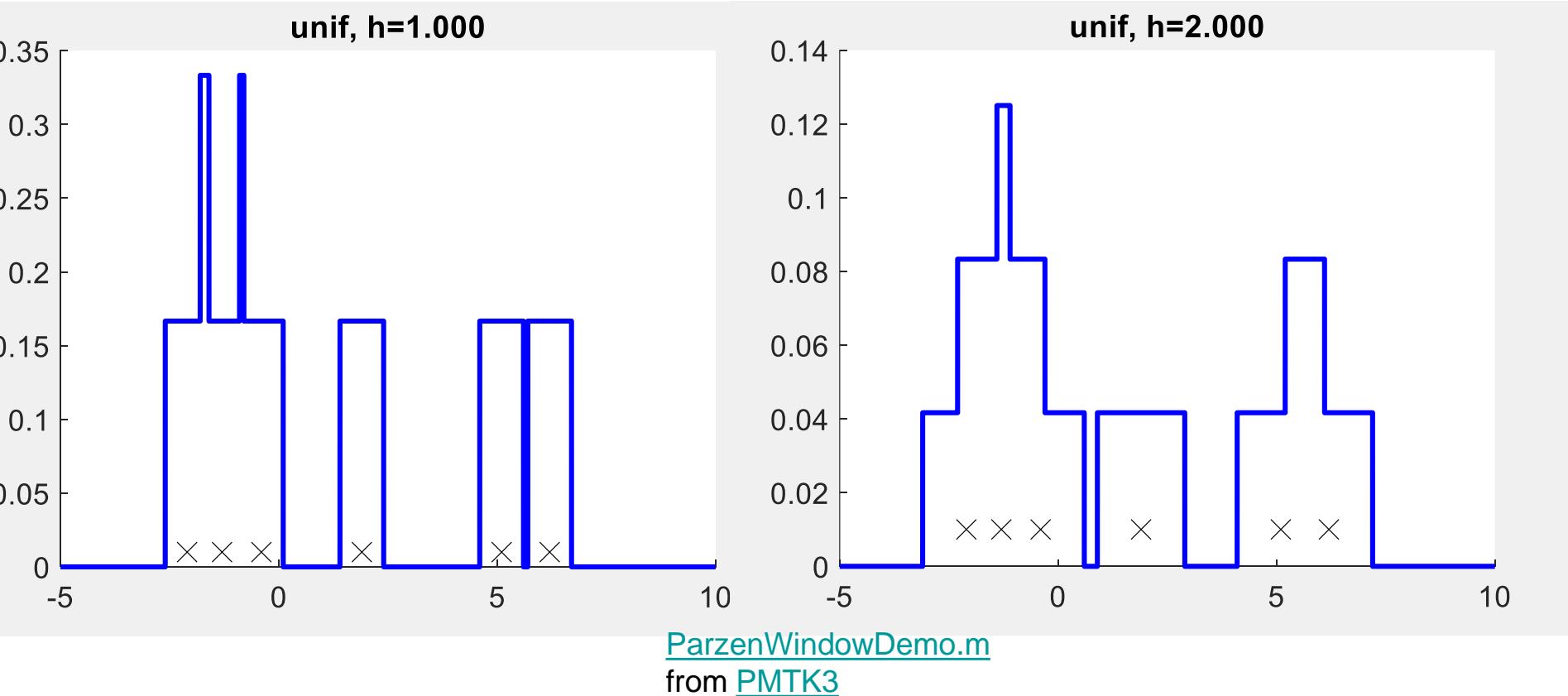
Kernel Density Estimation



[ParzenWindowDemo.m](#)
from [PMTK3](#)

Kernel Density Estimation

- ❑ A boxcar kernel results in a histogram estimation.
- ❑ Selection of h can be done with cross validation.



- Bowman, A. and A. Azzalini (1997). *Applied Smoothing Techniques for Data Analysis*. Oxford.

KDE for KNN Classifiers

- We use KDE for an alternative derivation of the KNN classifier.
- In KDE with a boxcar kernel, we fix the bandwidth and count how many data points fall within the hyper-cube centered on a datapoint.
- Suppose that we grow a volume around \mathbf{x} until we encounter K data points, regardless of their class label. Let the resulting volume have size $V(\mathbf{x})$, and let there be $N_c(\mathbf{x})$ examples from class c in this volume. We can estimate the class conditional density as follows:

$$p(\mathbf{x}|y = c, \mathcal{D}) = \frac{N_c(\mathbf{x})}{N_c V(\mathbf{x})}$$

where N_c is the total number of examples in class c in the whole data set. The class prior can be estimated by $p(y = c|\mathcal{D}) = \frac{N_c}{N}$

- Hence the class posterior is given by

$$p(y = c|\mathbf{x}, \mathcal{D}) = \frac{\frac{N_c(\mathbf{x})}{N_c V(\mathbf{x})} \frac{N_c}{N}}{\sum_{c'} \frac{N_{c'}(\mathbf{x})}{N_{c'} V(\mathbf{x})} \frac{N_{c'}}{N}} = \frac{N_c(\mathbf{x})}{\sum_{c'} N_{c'}(\mathbf{x})} = \frac{N_c(\mathbf{x})}{K}$$

Nadaraya-Watson Model

Interpolation with Noisy Inputs

- Consider the interpolation problem with noisy inputs rather than target values:

$$E = \frac{1}{2} \sum_{n=1}^N \int \left\{ y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n \right\}^2 \nu(\boldsymbol{\xi}) d\boldsymbol{\xi}$$

- We can optimize with respect to the function $y(\mathbf{x})$ to obtain:

$$y(\mathbf{x}) = \sum_{n=1}^N t_n h(\mathbf{x} - \mathbf{x}_n)$$

- The basis functions are given as (Nadaraya-Watson model):

$$h(\mathbf{x} - \mathbf{x}_n) = \frac{\nu(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^N \nu(\mathbf{x} - \mathbf{x}_n)}$$

- If noise distribution $\nu(\boldsymbol{\xi})$ is isotropic (function of $\|\boldsymbol{\xi}\|$), these basis functions would be radial.

- Bishop, C. M. (1995a). [Neural Networks for Pattern Recognition](#). Oxford University Press
- Webb, A. R. (1994). [Functional approximation by feed-forward networks: a least-squares approach to generalisation](#). *IEEE Transactions on Neural Networks* 5(3), 363–371.

Radial Basis Functions: Proof

- We take a variation of $y(\mathbf{x})$: $y(\mathbf{x}) \rightarrow y(\mathbf{x}) + \varepsilon\eta(\mathbf{x})$. Then:

$$E[y + \varepsilon\eta] = \frac{1}{2} \sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) + \varepsilon\eta(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\}^2 v(\boldsymbol{\xi}) d\boldsymbol{\xi}$$

- We expand in ε and set the coefficient in ε (functional 1st derivative) to zero:

$$\sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\} \eta(\mathbf{x}_n + \boldsymbol{\xi}) v(\boldsymbol{\xi}) d\boldsymbol{\xi} = 0$$

- Since this must be true for any variation $\eta(\mathbf{x})$, we choose:
 $\eta(\mathbf{x}) = \delta(\mathbf{x} - \mathbf{z})$. This leads to:

$$\sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\} \delta(\mathbf{x}_n + \boldsymbol{\xi} - \mathbf{z}) v(\boldsymbol{\xi}) d\boldsymbol{\xi} = \sum_{n=1}^N \{y(\mathbf{z}) - t_n\} v(\mathbf{z} - \mathbf{x}_n) = 0 \Rightarrow$$

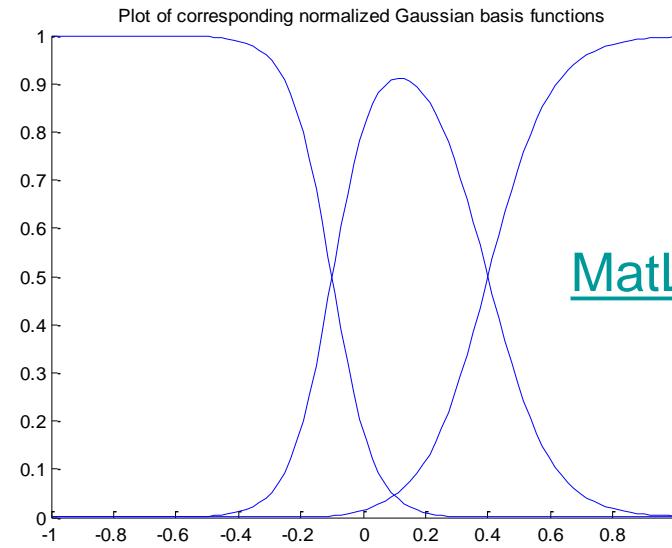
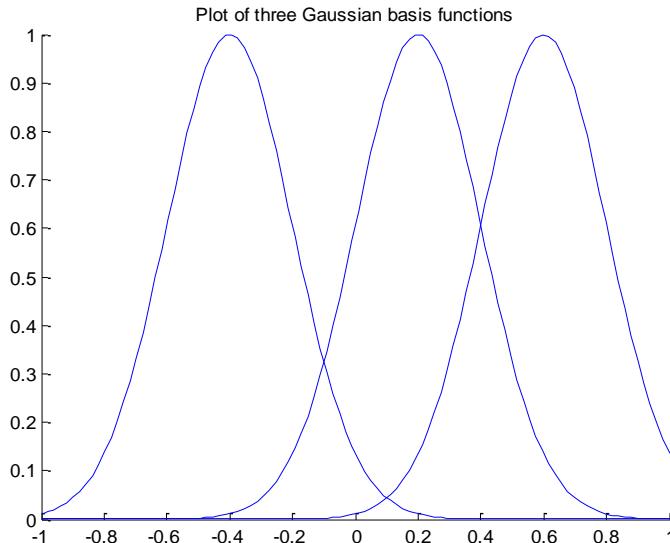
$$y(\mathbf{z}) = \sum_{n=1}^N t_n \left(v(\mathbf{z} - \mathbf{x}_n) / \sum_{l=1}^N v(\mathbf{z} - \mathbf{x}_l) \right)$$

Normalization Effect

- The functions in the Nadaraya-Watson model are normalized:

$$h(x - x_n) = \frac{\nu(x - x_n)}{\sum_{n=1}^N \nu(x - x_n)}, \quad \sum_n h(x - x_n) = 1$$

- It is useful to have normalized basis functions to avoid regions in an input space where all of the basis functions take small values (leading to small predictions or predictions controlled solely by the bias parameter).



[MatLab Code](#)

Radial Basis For Regression

- In the application of kernel density estimation to regression, there is one basis function associated with every data point, and the corresponding model can be computationally costly to evaluate when making predictions for new data points.
- Models have therefore been proposed using the expansion in radial basis functions but where the number M of basis functions is smaller than the number N of data points.
- Typically, the number of basis functions, and the centers μ_i are determined based on the input data $\{x_n\}$ alone. The basis functions are kept fixed and the coefficients $\{w_i\}$ are determined by least squares.

- Broomhead, D. S. and D. Lowe (1988). [Multivariable functional interpolation and adaptive networks](#). *Complex Systems* **2**, 321–355.
- Moody, J. and C. J. Darken (1989). [Fast learning in networks of locally-tuned processing units](#). *Neural Computation* **1**(2), 281–294.
- Poggio, T. and F. Girosi (1990). [Networks for approximation and learning](#). *Proceedings of the IEEE* **78**(9), 1481–1497.

Reducing the Size of the Basis

- To reduce the algorithmic cost, we can keep the number of basis functions M smaller than input data size N
- Centers locations μ_i are determined based on the input data $\{x_n\}$ alone
- Coefficients $\{w_i\}$ are determined by least squares
- For the choice of centers:
 - Random subset of the training data points
 - Orthogonal least squares (sequential) – the next data point to be chosen as a basis function center gives the greatest error reduction
 - Using clustering algorithms (K-means) – basis functions centers no longer coincide with the training data points.

▪ Chen, S., C. F. N. Cowan, and P. M. Grant (1991). [Orthogonal least squares learning algorithm for radial basis function networks](#). *IEEE Transactions on Neural Networks* 2(2), 302–309.

Parzen Density Estimator

- We have seen in an earlier lecture the prediction of the linear regression model - linear combination of t_n with 'equivalent kernel' values
- The same result is obtained from the Parzen density estimator, $f(\mathbf{x}, t)$ is the component density function centered at each training data point

$$p(\mathbf{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x} - \mathbf{x}_n, t - t_n)$$

- Regression function

$$\begin{aligned} y(\mathbf{x}) &= \mathbb{E}[t | \mathbf{x}] = \int_{-\infty}^{+\infty} tp(t | \mathbf{x}) dt = \\ &= \frac{\int tp(\mathbf{x}, t) dt}{\int p(\mathbf{x}, t) dt} = \frac{\sum_n \int tf(\mathbf{x} - \mathbf{x}_n, t - t_n) dt}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt} \end{aligned}$$

Kernel Regression: Nadaraya-Watson Model

- Assume that the component density functions have zero mean so that

$$\int_{-\infty}^{+\infty} tf(x, t) dt = 0$$

for all values of x . Then by variable change

$$y(x) = \frac{\sum_n \int (t - t_n + t_n) f(x - x_n, t - t_n) dt}{\sum_m \int f(x - x_m, t - t_m) dt} = \frac{\sum_n g(x - x_n) t_n}{\sum_m g(x - x_m)} = \\ = \sum k(x, x_n) t_n$$

with $g(x) = \int_{-\infty}^{+\infty} f(x, t) dt$ and $k(x, x_n)$ given by

$$k(x, x_n) = \frac{g(x - x_n)}{\sum_m g(x - x_m)}$$

Kernel Regression: Nadaraya-Watson Model

- For localized kernels, this regression model gives more weight to the data points \mathbf{x}_n close to \mathbf{x} .

$$y(\mathbf{x}) = \sum_n k(\mathbf{x}, \mathbf{x}_n) t_n$$

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)}$$

- The kernel satisfies the summation constraint:

$$\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1$$

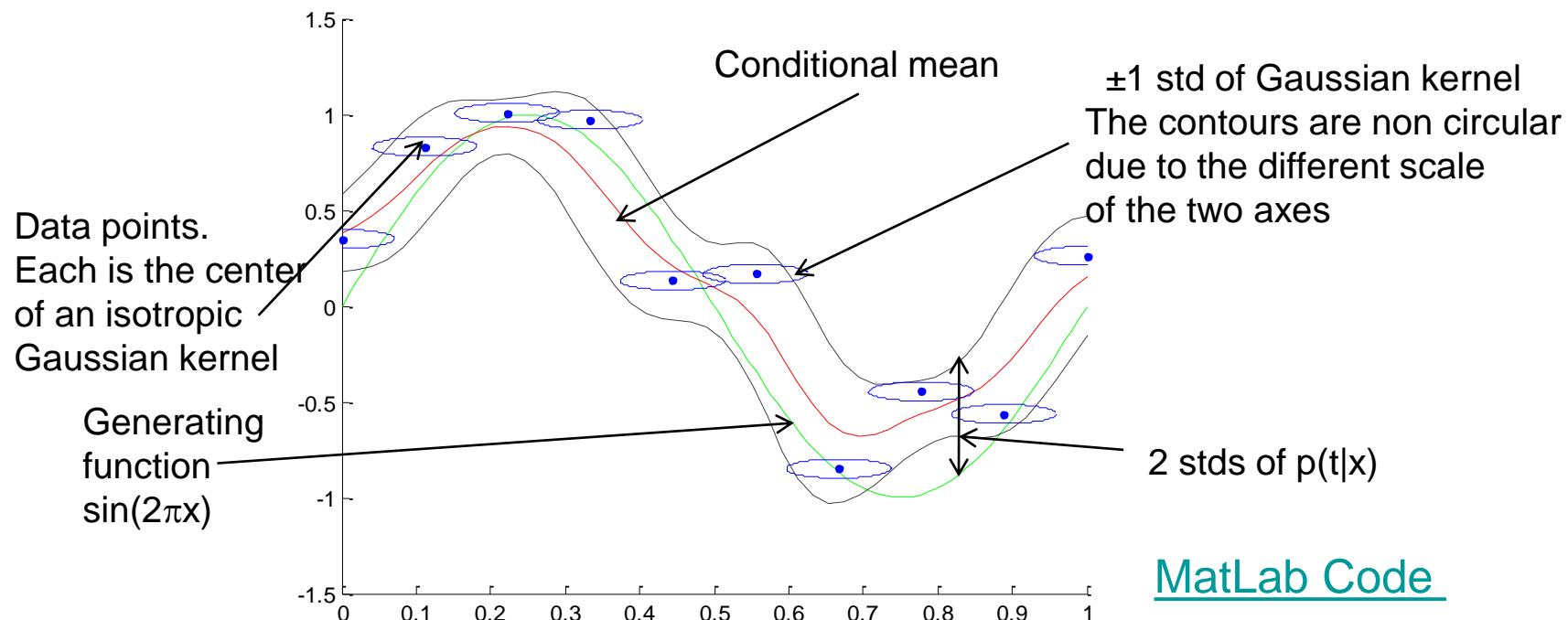
- In addition to the conditional expectation $y(\mathbf{x})$, the model provides the full conditional distribution:

$$p(t | \mathbf{x}) = \frac{p(t, \mathbf{x})}{\int p(t, \mathbf{x}) dt} = \frac{\sum_n f(\mathbf{x} - \mathbf{x}_n, t - t_n)}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt}$$

- For Gaussian component distributions, this is a mixture of Gaussians.

Kernel Regression: Nadaraya-Watson Model

- Single input variable x . $f(x, t)$ is a zero-mean isotropic Gaussian kernels over the variable $z = (x, t)$ with variance σ^2



- Nadaraya, E. A. (1964). [On estimating regression](#). *Theory of Probability and its Applications* **9**(1), 141–142.
- Watson, G. S. (1964). [Smooth regression analysis](#). *Sankhya: The Indian Journal of Statistics*. Series A **26**, 359–372.

Kernel Regression

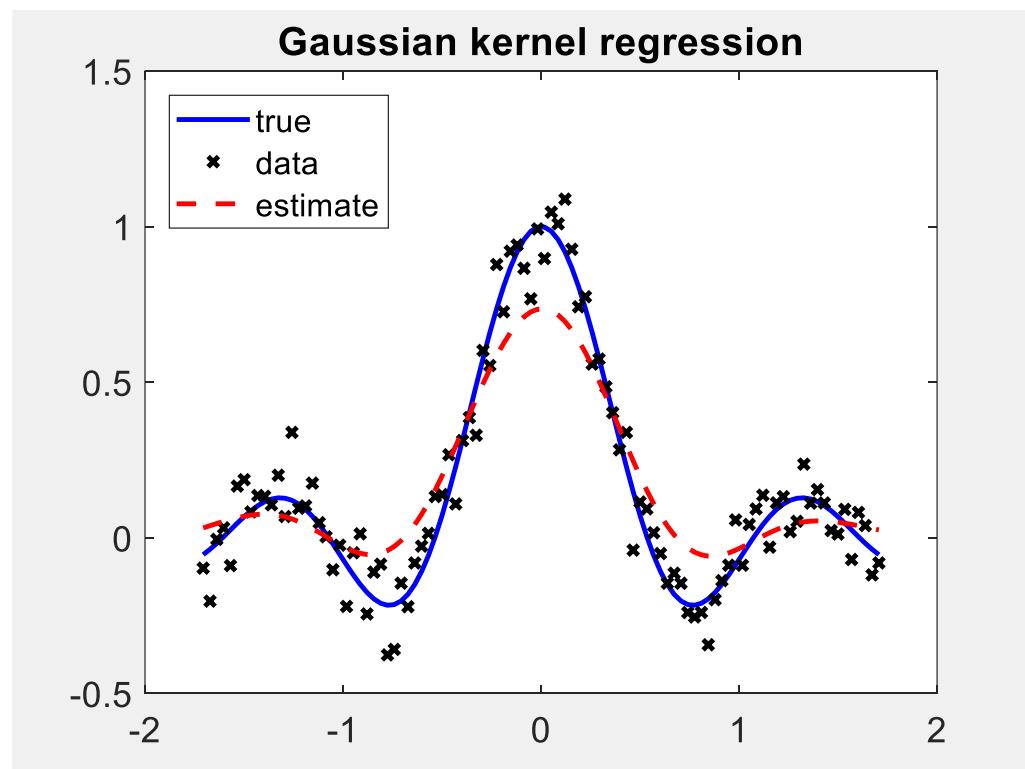
- This method only has one free parameter, namely h . For 1d data, if the true density is Gaussian and we are using Gaussian kernels, the optimal h is given by

$$h = \left(\frac{4}{3N}\right)^{1/5} \hat{\sigma}, \hat{\sigma} = 1.4826 \text{ MAD}, \text{MAD} = \text{median}(|x - \text{median}(x)|)$$

- We also take

$$h_x = h_t (= \sqrt{h_x h_t})$$

[kernelRegressionDemo.m](#)
from [PMTK3](#)



Locally Weighted Regression

- If we define $k_h(x - x_i) = k(x, x_i)$, we can rewrite the prediction made by kernel regression as

$$\hat{f}(x_*) = \sum_{i=1}^N y_i \frac{k(x_*, x_i)}{\sum_{i'=1}^N k(x_*, x_{i'})}$$

- If $k(x, x_i)$ is not smoothing kernel we can write:

$$\hat{f}(x_*) = \sum_{i=1}^N y_i k(x_*, x_i)$$

- This model is essentially fitting a constant function locally. We can fit a linear regression model for each x_* by solving

$$\max_{\beta(x_*)} \sum_{i=1}^N k(x_*, x_i)(y_i - \beta(x_*)^T \phi(x_i))^2, \quad \phi(x_i) = [1, x]$$

- This leads to a weighted LS problem $\beta(x_*) = (\Phi^T D(x_*) \Phi)^{-1} \Phi^T D(x_*) y$
 $D = diag(k(x_*, x_i))$, Φ is an $N \times (D + 1)$ design matrix
- Predictions take the form:

$$\hat{f}(x_*) = \phi(x_*)^T \beta(x_*) = \phi(x_*)^T (\Phi^T D(x_*) \Phi)^{-1} \Phi^T D(x_*) y = \sum_{i=1}^N w_i(x_*) y_i$$

- Cleveland, W. and S. Devlin (1988). [Locally-weighted regression: An approach to regression analysis by local fitting](#). *J. of the Am. Stat. Assoc.* 83(403), 596–610.
- Edakunni, N., S. Schaal, and S. Vijayakumar (2010). [Probabilistic incremental locally weighted learning using randomly varying coefficient model](#). Technical report, USC.

Kernel Regression: Nadaraya-Watson Model

- We can use more flexible forms of Gaussian components, e.g. having different variance parameters for the input and target variables.
 - We could alternatively model $p(t, x)$ using a Gaussian mixture model, trained using EM techniques and then find $p(t|x)$.
 - ✓ In this case we don't have a representation in terms of kernel functions evaluated at the training set data points.
 - ✓ However, the number of components in the mixture model can be smaller than the number of training set points, resulting in a model that is faster to evaluate for test data points.
-
- [Ghahramani, Z. and M. I. Jordan \(1994\). Supervised learning from incomplete data via an EM approach.](#) In J. D. Cowan, G. T. Tesauro, and J. Alspector (Eds.), [Advances in Neural Information Processing Systems, Volume 6](#), pp. 120–127.

The Kernel Trick

- **Kernel Trick:** We can work with the feature vectors x , but modify the algorithm so that it replaces all inner products of the form $\langle x, x' \rangle$ with a call to the (Mercer) kernel function, $k(x, x')$.
- **Kernelized nearest neighbor classification:** Here we need to compute the Euclidean distance of a test vector to all the training points, find the closest one, and look up its label. This can be kernelized by observing that

$$\|x_i - x_{i'}\| = \langle x_i, x_i \rangle + \langle x_{i'}, x_{i'} \rangle - 2\langle x_i, x_{i'} \rangle$$

- This allows us to apply the nearest neighbor classifier to structured data objects.

Kernelized K-Medoids Clustering

- The K-medoids algorithm is similar to K-means - instead of representing each cluster's centroid by the mean of all data vectors assigned to the cluster, we make each centroid to be one of the data vectors themselves.
- Thus we always deal with integer indexes, rather than data objects.
- We assign objects to their closest centroids as before. When we update the centroids, we look at each point that belongs to the cluster, and measure the sum of its distances to all the others in the same cluster; we then pick the one which has the smallest such sum:

$$m_k = \operatorname{argmin}_{i:z_i=k} \sum_{i':z_{i'}=k} d(i, i'), \quad d(i, i') = \|x_i - x_{i'}\|_2^2$$

Kernelized K-Medoids Clustering

- This takes $\mathcal{O}(n_k^2)$ work per cluster, whereas K-means takes $\mathcal{O}(n_k D)$ to update each cluster.
- This method can be modified to derive a classifier, by computing the nearest medoid for each class. This is known as **nearest medoid classification**. It can be kernelized by using

$$\|x_i - x_{i'}\| = \langle x_i, x_i \rangle + \langle x_{i'}, x_{i'} \rangle - 2\langle x_i, x_{i'} \rangle$$

- **1** initialize $m_{1:K}$ as a random subset of size K from $\{1, \dots, N\}$;
- **2 repeat**
- **3** $z_i = \operatorname{argmin}_k d(i, m_k)$ for $i = 1 : N$;
- **4** $m_k \leftarrow \operatorname{argmin}_{i:z_i=k} \sum_{i':z_{i'}=k} d(i, i')$, for $k = 1 : K$;
- **5 until converged;**
 - Hastie, T., R. Tibshirani, and J. Friedman (2009). *The Elements of Statistical Learning*. Springer. 2nd edition

Kernel PCA

Kernel PCA

- Kernel trick: take an algorithm expressed in terms of scalar products of the form $x^T x$ and generalize it by replacing the scalar products with a nonlinear kernel.
- Apply this technique of kernel substitution to PCA, obtaining a nonlinear algorithm called kernel PCA ([Scholkopf et al. 1998](#))
- Consider a data set $\{x_n\}, n = 1, \dots, N$, of dimension D . Assume that we have already subtracted the sample mean from each x_n , so that $\sum_n x_n = 0$.
- Start by expressing PCA in such a form that the data $\{x_n\}$ appear only in the form of the scalar products $x_n^T x_m$.
 - Scholkopf, B., A. Smola, and K.-R. Muller (1998). [Nonlinear component analysis as a kernel eigenvalue problem](#). *Neural Computation* 10(5), 1299–1319.

Kernel PCA

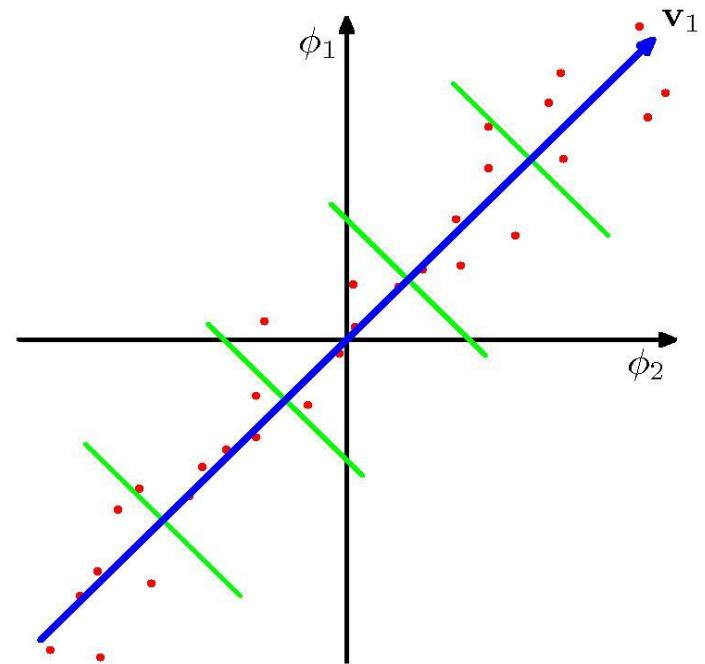
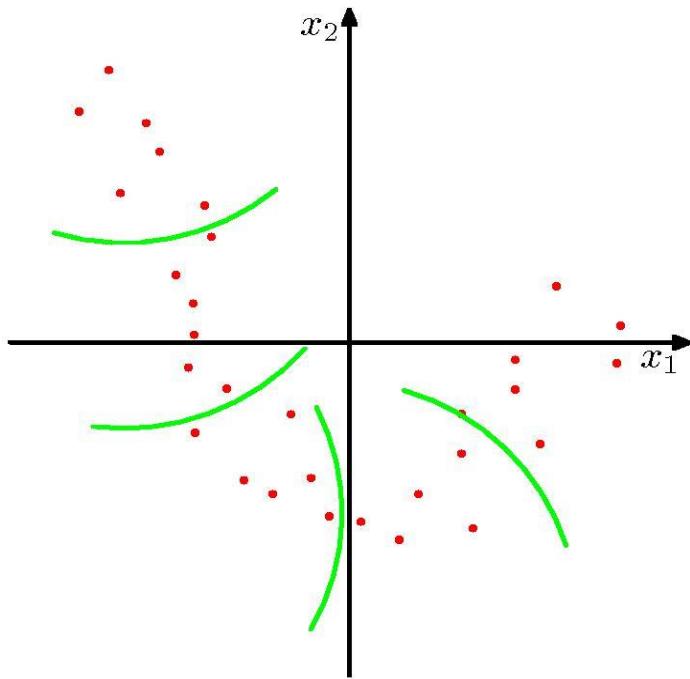
- Recall that the principal components are defined by the normalized eigenvectors \mathbf{u}_i of the $N \times N$ covariance matrix

$$S\mathbf{u}_i = \lambda_i \mathbf{u}_i, \text{ where } S = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

- Consider a nonlinear transformation $\phi(\mathbf{x})$ into an M -dimensional feature space, so that each data point \mathbf{x}_n is thereby projected onto a point $\phi(\mathbf{x}_n)$.
- Now perform standard PCA in the feature space.
 - This implicitly defines a nonlinear principal component model in the original data space.

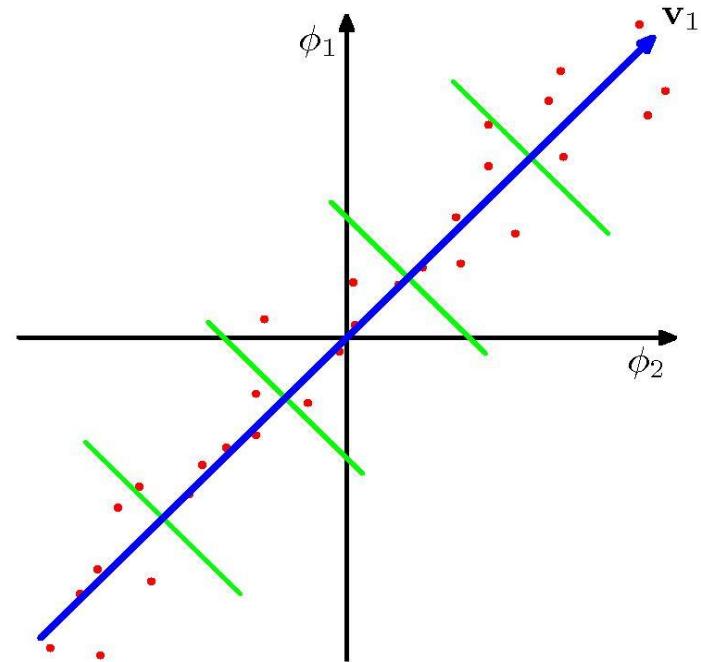
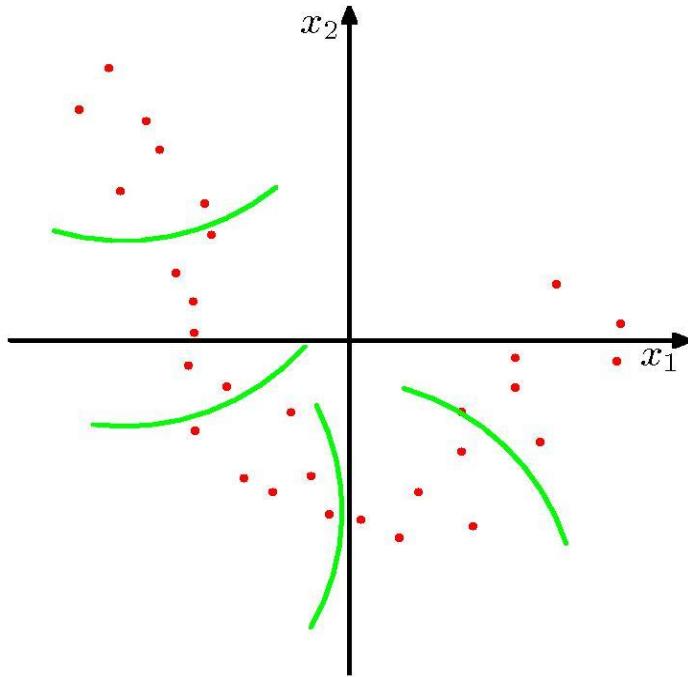
Kernel PCA

- A data set in the original data space (left) is projected by $\phi(x)$ into a feature space (right).
- By performing PCA in the feature space, we obtain the principal components (v_1 being the 1st one).



Kernel PCA

- The green lines in feature space indicate the linear projections onto v_1 , which correspond to nonlinear projections in the original data space.
- In general, it is not possible to represent the nonlinear principal component by a vector in x space.



kPCA

- Recall how we could compute a low-dimensional linear embedding of some data using PCA.
- This required finding the eigenvectors of the sample covariance matrix $S = \frac{1}{N} X^T X$.
- However, we can also compute PCA by finding the eigenvectors of the inner product matrix XX^T , as we show below. This will allow us to produce a nonlinear embedding, using the kernel trick.
- First, let U be an orthogonal matrix containing the eigenvectors of XX^T with corresponding eigenvalues in Λ . By definition we have $(XX^T)U = U\Lambda$. Pre-multiplying by X^T gives

$$(X^T X)(X^T U) = (X^T U)\Lambda$$

N^{*}N

from which we see that the eigenvectors of $X^T X$ (and hence of S) are $V = X^T U$, with eigenvalues given by Λ as before. However, these eigenvectors are not normalized, since

$$\|v_j\|^2 = u_j^T X X^T u_j = \lambda_j u_j^T u_j = \lambda_j.$$

- So the normalized eigenvectors are given by $V_{pca} = X^T U \Lambda^{-1/2}$.

kPCA

- This is a useful trick for PCA if $D > N$, since $\mathbf{X}^T \mathbf{X}$ has size $D \times D$, whereas $\mathbf{X} \mathbf{X}^T$ has size $N \times N$. It also allow us to use the kernel trick.
- Let $\mathbf{K} = \mathbf{X} \mathbf{X}^T$ be the Gram matrix. From Mercer's theorem the use of a kernel implies some underlying feature space, so we are implicitly replacing x_i with $\phi(x_i) = \phi_i$. Let Φ be the corresponding design matrix, and $S_\phi = \frac{1}{N} \sum_i \phi_i \phi_i^T$ be the corresponding covariance matrix in feature space. The eigenvectors are given by $\mathbf{V}_{kPCA} = \Phi^T \mathbf{U} \Lambda^{-1/2}$, where \mathbf{U}, Λ contain the eigenvectors & eigenvalues of \mathbf{K} .
- We cannot compute \mathbf{V}_{kPCA} , since ϕ_i is potentially infinite dimensional. However, we can compute the projection of a test vector x_* onto the feature space as follows:

$$\phi_*^T \mathbf{V}_{kPCA} = \phi_*^T \Phi^T \mathbf{U} \Lambda^{-1/2} = \mathbf{k}_*^T \mathbf{U} \Lambda^{-1/2}$$

where $\mathbf{k}_* = [\kappa(x_*, x_1), \dots, \kappa(x_*, x_N)]$.

- So far, we have assumed the projected data has zero mean, which is not the case in general.

Kernel PCA

- For the moment, let us assume that the projected data set also has zero mean, so that $\sum_{n=1}^N \phi(x_n) = 0$. We shall return to this point shortly.
- The $D_* \times D_*$ sample covariance matrix in feature space is given by

$$C = \frac{1}{N} \sum_{n=1}^N \phi(x_n) \phi(x_n)^T$$

- The eigenvector expansion is defined as:

$$C\mathbf{v}_i = \lambda_i \mathbf{v}_i, i = 1, \dots, D_*$$

- Can we solve this eigenvalue problem without having to work explicitly in the feature space?

Kernel PCA

$$C = \frac{1}{N} \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \quad C\mathbf{v}_i = \lambda_i \mathbf{v}_i, i = 1, \dots, D_*$$

- From the definition of C , the eigenvector equations tells us that \mathbf{v}_i satisfies

$$\frac{1}{N} \sum_{n=1}^N \phi(x_n) \left\{ \phi(x_n)^T \mathbf{v}_i \right\} = \lambda_i \mathbf{v}_i, i = 1, \dots, D_*$$

- We see that (provided $\lambda_i > 0$), the vector \mathbf{v}_i is given by a linear combination of the $\phi(x_n)$ and so can be written in the form

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(x_n), i = 1, \dots, D_*$$

Kernel PCA

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \left\{ \phi(\mathbf{x}_n)^T \mathbf{v}_i \right\} = \lambda_i \mathbf{v}_i, i = 1, \dots, D_*$$

- Substituting

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n), i = 1, \dots, D_*$$

back into the eigenvector equation, we obtain

$$\left[\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right] \sum_{m=1}^N a_{im} \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n), i = 1, \dots, D_*$$

- The key step is now to express this in terms of the kernel function

$$K_{nm} \equiv k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

which we do by multiplying both sides by $\phi(\mathbf{x}_l)^T$.

Kernel PCA

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \sum_{m=1}^N a_{im} \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n), i = 1, \dots, D_*$$

- Multiplying both sides by $\phi(\mathbf{x}_l)^T$ to give

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_l)^T \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \sum_{m=1}^N a_{im} \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} \phi(\mathbf{x}_l)^T \phi(\mathbf{x}_n), i = 1, \dots, D_* \Rightarrow$$

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_l, \mathbf{x}_n) \sum_{m=1}^N a_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} k(\mathbf{x}_l, \mathbf{x}_n), i = 1, \dots, D_* \Rightarrow$$

- This can be written in a matrix form as:

$$K^2 \mathbf{a}_i = \lambda_i N K \mathbf{a}_i$$

\mathbf{a}_i N -dimensional column vector with elements $a_{in}, n = 1, \dots, N$.

- Can simplify $K \mathbf{a}_i = \lambda_i N \mathbf{a}_i$. These Eqs. differ by eigen-vectors of K with 0 eigenvalues that don't affect the principal components projection.

Kernel PCA

$$\mathbf{K}\mathbf{a}_i = \lambda_i N\mathbf{a}_i$$

- We can find solutions for \mathbf{a}_i by solving the above eigenvalue problem.
- The normalization condition for the coefficients \mathbf{a}_i is obtained by requiring that the eigenvectors in feature space be normalized.
- Using

and

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n), i = 1, \dots, D_*$$

we have:

$$\mathbf{K}\mathbf{a}_i = \lambda_i N\mathbf{a}_i$$

$$1 = \mathbf{v}_i^T \mathbf{v}_i = \frac{1}{N} \sum_{n=1}^N \sum_{m=1}^N a_{in} a_{im} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = \mathbf{a}_i^T \mathbf{K} \mathbf{a}_i = \lambda_i N \mathbf{a}_i^T \mathbf{a}_i \Rightarrow 1 = \lambda_i N \mathbf{a}_i^T \mathbf{a}_i$$

Kernel PCA

- The principal component projections can now be cast in terms of the kernel function.
- Using

$$\boldsymbol{v}_i = \sum_{n=1}^N a_{in} \boldsymbol{\phi}(\boldsymbol{x}_n), i = 1, \dots, D_*$$

the projection of a point \boldsymbol{x} onto eigenvector i is given by

$$y_i(\boldsymbol{x}) = \boldsymbol{\phi}(\boldsymbol{x})^T \boldsymbol{v}_i = \sum_{n=1}^N a_{in} \boldsymbol{\phi}(\boldsymbol{x})^T \boldsymbol{\phi}(\boldsymbol{x}_n) = \sum_{n=1}^N a_{in} k(\boldsymbol{x}, \boldsymbol{x}_n)$$

Kernel PCA

- In the original D -dimensional x space: D orthogonal eigenvectors and thus at most D linear principal components.
- The dimensionality D_* of the feature space, however, can be much larger than D (even infinite), and thus we can find a number of nonlinear principal components that can exceed D .
- Note, however, that the number of nonzero eigenvalues cannot exceed the number N of data points:
 - The covariance matrix in feature space has rank at most equal to N .
 - Recall that kernel PCA involves the eigenvector expansion of the $N \times N$ matrix K .

Kernel PCA

- So far we have assumed that the projected data set given by $\phi(x_n)$ has zero mean, which in general will not be the case.
- We cannot simply compute and then subtract off the mean, since we want to formulate the algorithm purely in terms of the kernel function.
- The projected data points after centralizing are given by

$$\phi(x_n) = \phi(x_n) - \frac{1}{N} \sum_{l=1}^N \phi(x_l)$$

- The corresponding elements of the Gram matrix $\tilde{K}_{nm} = \tilde{\phi}(x_n)^T \tilde{\phi}(x_m)$ can be computed as shown next.

Kernel PCA

$$\begin{aligned}\tilde{K}_{nm} &= \tilde{\phi}(x_n)^T \tilde{\phi}(x_m) = \left(\phi(x_n) - \frac{1}{N} \sum_{l=1}^N \phi(x_l) \right)^T \left(\phi(x_m) - \frac{1}{N} \sum_{l=1}^N \phi(x_l) \right) \\ &= \phi(x_n)^T \phi(x_m) - \frac{1}{N} \sum_{l=1}^N \phi(x_n)^T \phi(x_l) \\ &\quad - \frac{1}{N} \sum_{l=1}^N \phi(x_l)^T \phi(x_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N \phi(x_j)^T \phi(x_l) \\ &= k(x_n, x_m) - \frac{1}{N} \sum_{l=1}^N k(x_n, x_l) - \frac{1}{N} \sum_{l=1}^N k(x_l, x_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N k(x_j, x_l)\end{aligned}$$

□ We can express this in matrix notation as:

$$\tilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N$$

where $\mathbf{1}_N$ denotes the $N \times N$ matrix in which every element takes the value $1/N$.

Kernel PCA

$$\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$$

- Thus we can evaluate $\tilde{\mathbf{K}}$ using only the kernel function and then use $\tilde{\mathbf{K}}$ to determine the eigenvalues and eigenvectors.
- Note that the standard PCA algorithm is recovered as a special case if we use a linear kernel $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$.
- We often use Gaussian kernels of the form:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|}{2\sigma^2}\right)$$

- Whereas linear PCA is limited to using $M \leq D$ components, in kPCA, we can use up to N components, since Φ is $N \times D_*$, where D_* is the (potentially infinite) dimensionality of the embedded feature vectors.

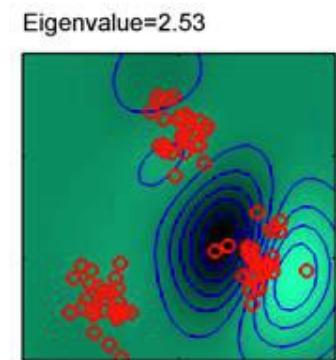
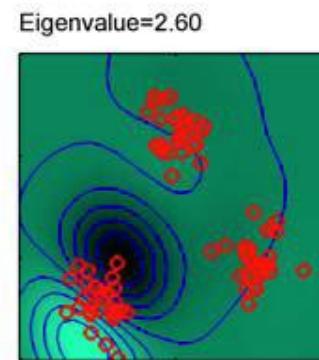
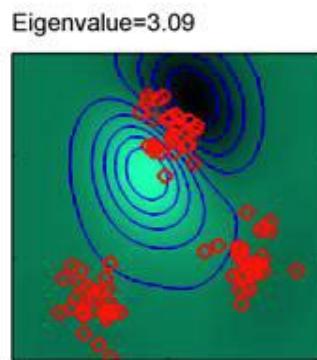
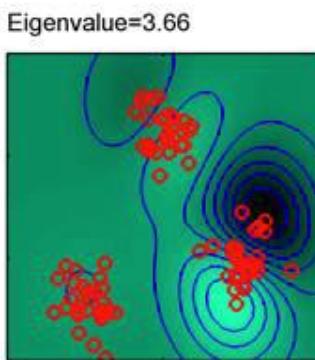
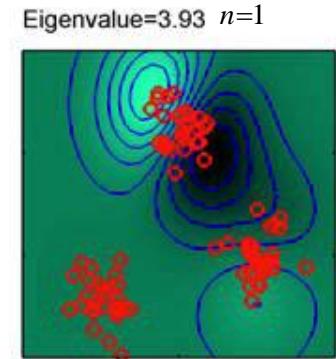
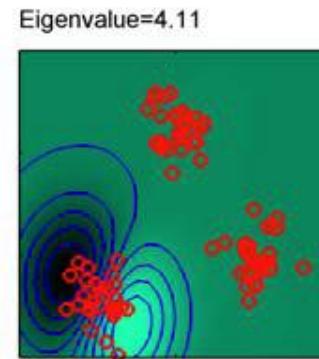
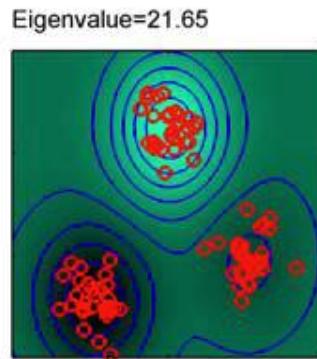
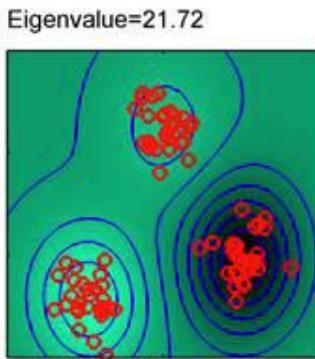
kPCA - Algorithm

- 1 Input: K of size $N \times N$, K_* of size $N_* \times N$, number of latent dimensions M ;
- 2 $O = \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$;
- 3 $\tilde{K} = K - OK - KO + OKO$;
- 4 $[U, \Lambda] = eig(\tilde{K})$;
- 5 **for** $i = 1 : N$ **do**
- 6 $v_i = u_i / \sqrt{\lambda_i}$;
- 7 $O_* = \frac{1}{N} \mathbf{1}_{N_*N} \mathbf{1}_{N_*}^T$;
- 8 $\tilde{K}_* = K_* - O_* K_* - K_* O_* + O_* K_* O_*$;
- 9 $Z = \tilde{K}_* V(:, 1 : M)$

Kernel PCA

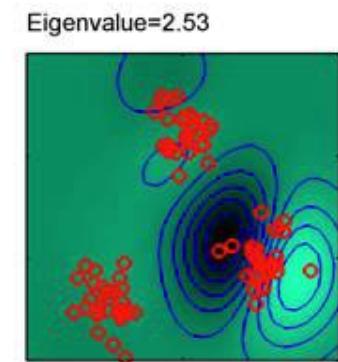
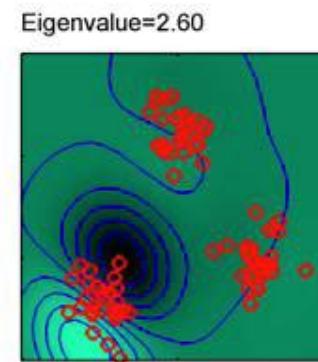
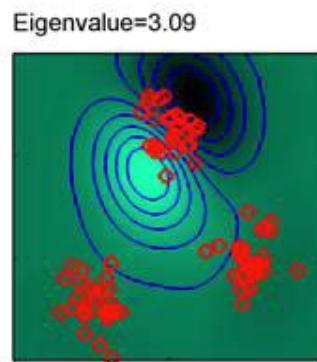
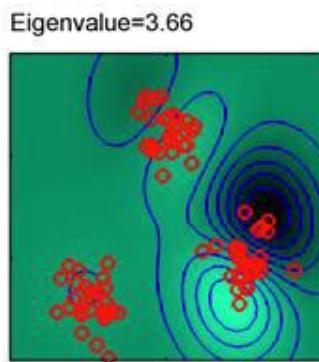
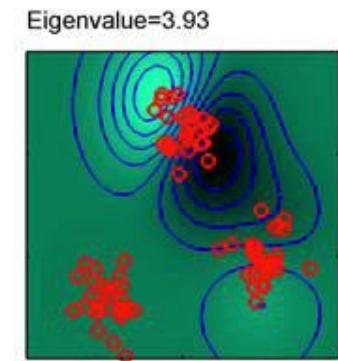
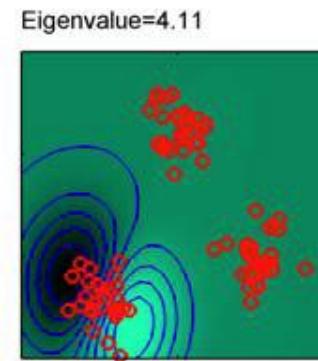
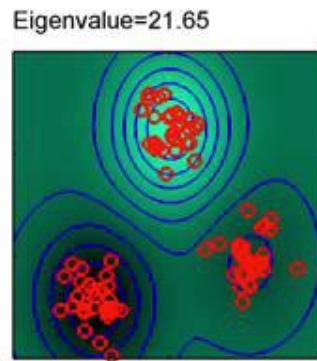
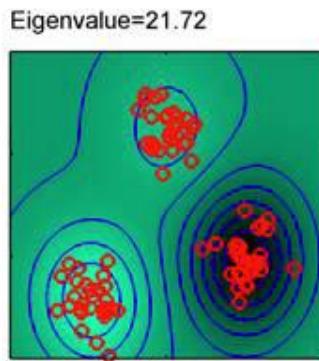
- ❑ KPCA, with a Gaussian kernel ($2\sigma^2 = 0.1$) applied to a synthetic data in 2D, shows the first 8 eigenfunctions and eigenvalues.
- ❑ The lines along which the projection onto the corresponding principal component is constant are shown:

$$\phi(\mathbf{x})^T \mathbf{v}_i = \sum_{n=1}^N a_{in} k(\mathbf{x}, \mathbf{x}_n)$$



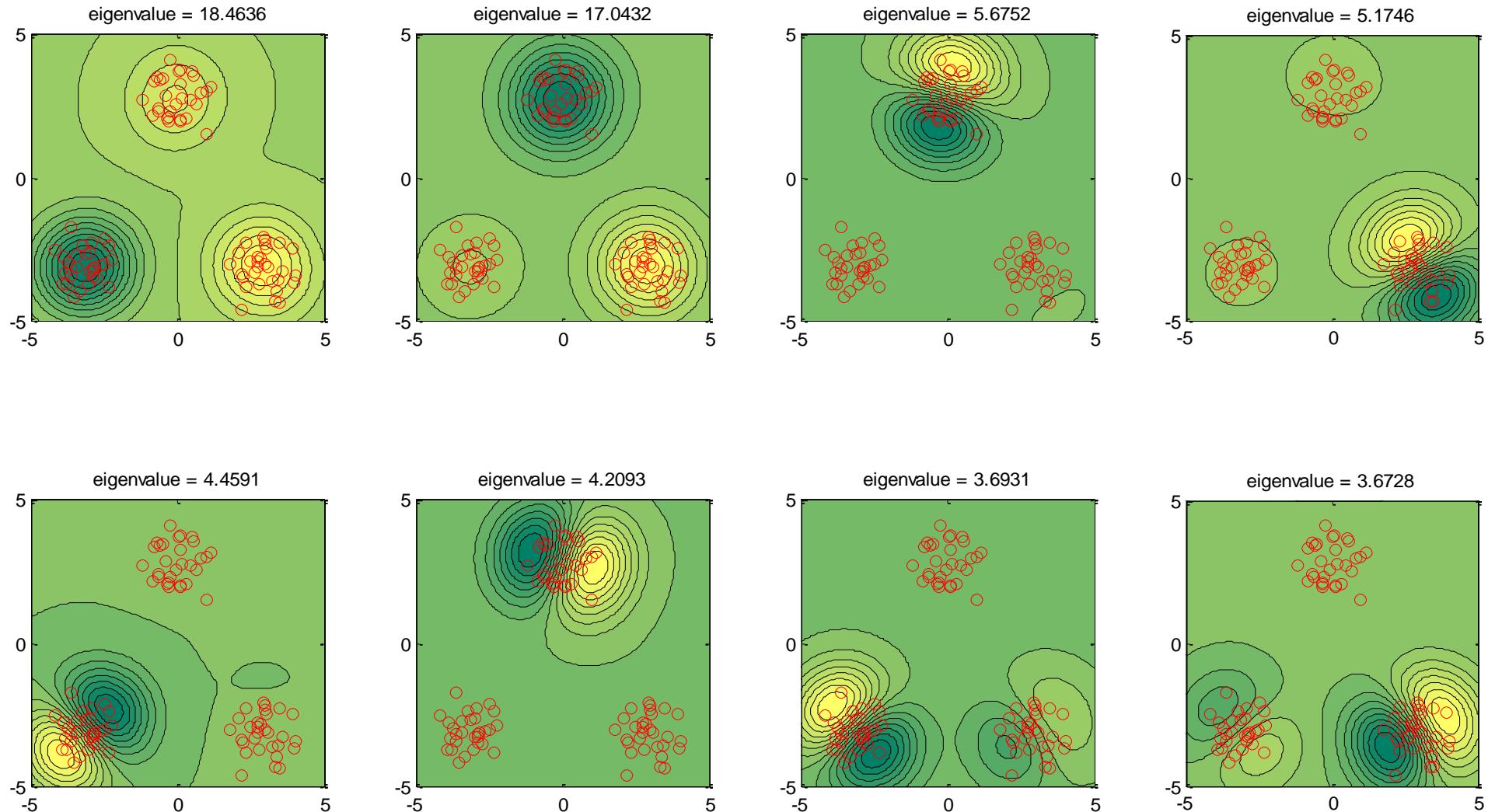
Kernel PCA

- ❑ Note the first 2 eigenvectors separate the 3 clusters, the next three eigenvectors split each of the cluster into halves, and the following three eigenvectors again split the clusters into halves along directions orthogonal to the previous splits.



Matlab Implementation

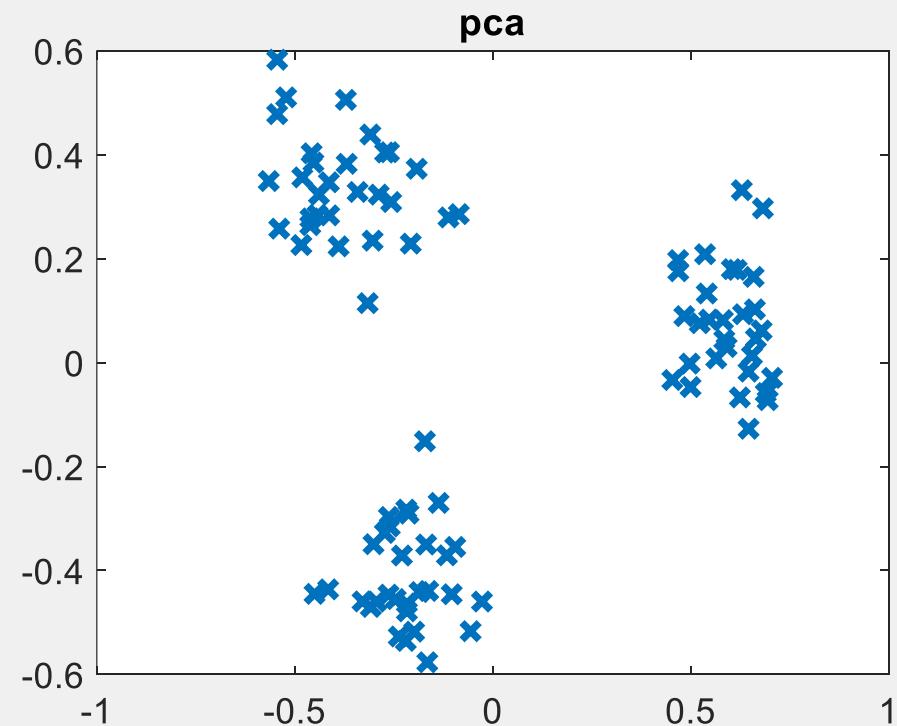
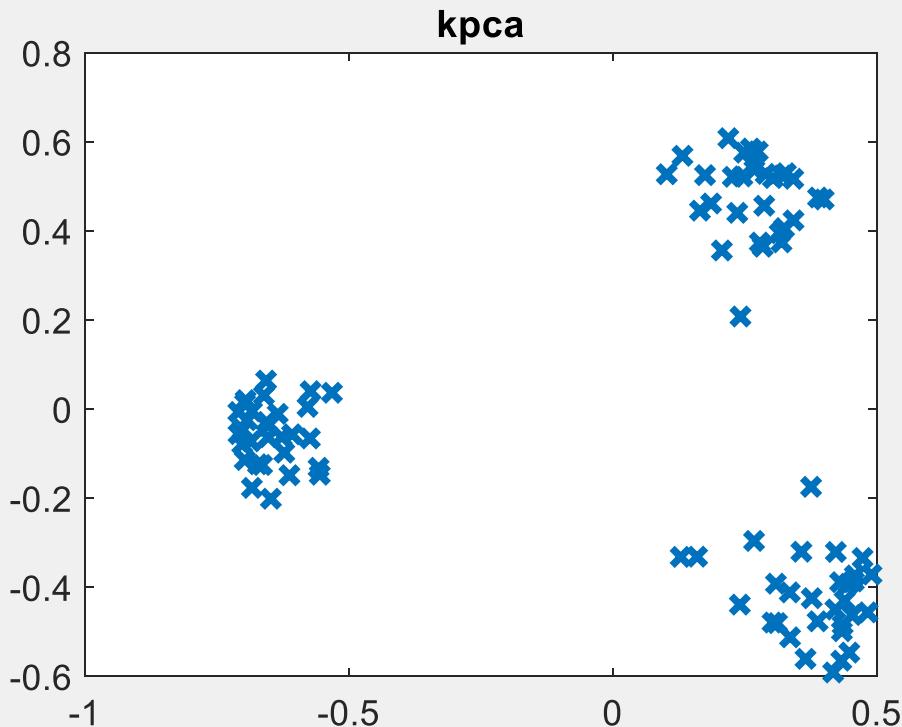
Kernel PCA



Matlab Implementation

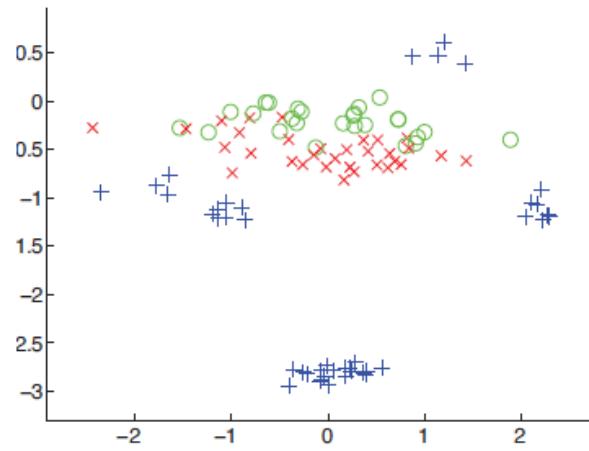
Kernel PCA

- Projection of the data from onto the first 2 principal bases computed using PCA and kPCA. PCA perfectly represents the data. kPCA represents each cluster by a different line.

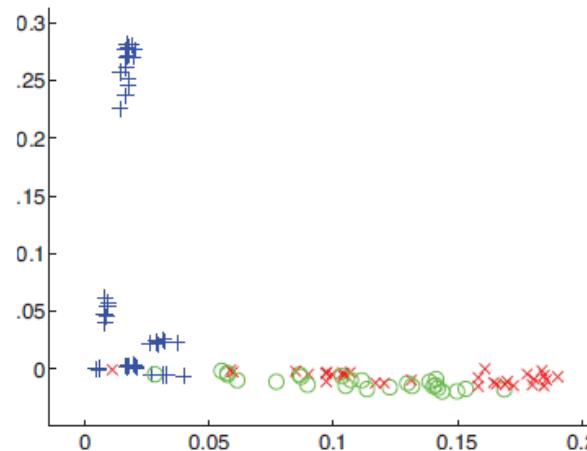


Oil Data Set

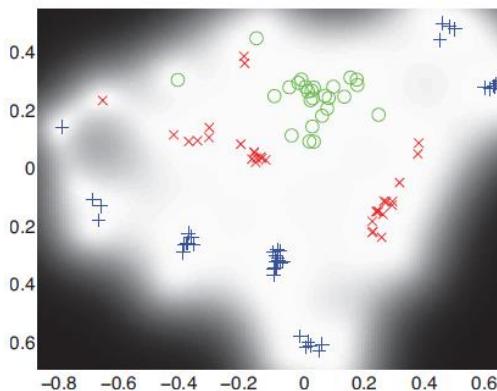
- We will use a 12 dimensional data set representing the three known phases of flow in an oil pipeline.
- We project this into 2d using PCA and kPCA (with an RBF kernel).
- If we perform nearest neighbor classification in the low-dimensional space, kPCA makes 13 errors and PCA makes 20.



PCA



kPCA



GPLVM with
Gaussian kernel

- Lawrence, N. D. (2005). [Probabilistic non-linear principal component analysis with gaussian process latent variable models](#). *J. of Machine Learning Research* 6, 1783–1816.

Kernel PCA: Disadvantages

- KPCA involves finding the eigenvectors of the $N \times N$ matrix \tilde{K} rather than the $D \times D$ matrix S of PCA.
 - In practice, for large N approximations are used.
- In PCA, we retain some reduced number $M < D$ of eigenvectors and then approximate a data vector x_n by its projection \hat{x}_n onto the M -dimensional principal subspace, defined by

$$\hat{x}_n = \sum_{i=1}^M (x_n^T u_i) u_i$$

- In KPCA, this in general is not possible.

Kernel PCA: Preimage problem

$$\hat{x}_n = \sum_{i=1}^M (x_n^T u_i) u_i$$

- This is not possible in KPCA. Indeed note that the mapping $\phi(x)$ maps the D -dimensional x space into a D -dimensional manifold in the D_* -dimensional feature space ϕ .
- The vector x is known as the pre-image of $\phi(x)$. The projection of points in feature space onto the linear PCA subspace in that space will typically not lie on the nonlinear D -dimensional manifold and so will not have a corresponding pre-image in data space.
- Techniques have been proposed for finding approximate pre-images (Bakir et al., 2004)
 - Bakir, G. H., J. Weston, and B. Scholkopf (2004). Learning to find pre-images. In S. Thrun, L. K. Saul, and B. Scholkopf (Eds.), Advances in Neural Information Processing Systems, Volume 16, pp. 449–456. MIT Press.